

Insights into Quasar UV Spectra Using Unsupervised Clustering Analysis

A. Tammour^{★1}, S. C. Gallagher^{1,2}, M. Daley³, & G. T. Richards⁴

¹ *Department of Physics and Astronomy, University of Western Ontario, 1151 Richmond St, London, ON, N6A 3K7, Canada*

² *Centre for Planetary and Space Exploration, University of Western Ontario, 1151 Richmond St, London, ON, N6A 3K7, Canada*

³ *Department of Computer Science, University of Western Ontario, 1151 Richmond St, London, ON, N6A 3K7, Canada*

⁴ *Department of Physics, Drexel University, Philadelphia, PA 19104, USA*

11 March 2016

ABSTRACT

Machine learning techniques can provide powerful tools to detect patterns in multi-dimensional parameter space. We use K-means –a simple yet powerful unsupervised clustering algorithm which picks out structure in unlabeled data– to study a sample of quasar UV spectra from the Quasar Catalog of the 10th Data Release of the Sloan Digital Sky Survey (SDSS-DR10) of [Pâris et al. \(2014\)](#). Detecting patterns in large datasets helps us gain insights into the physical conditions and processes giving rise to the observed properties of quasars. We use K-means to find clusters in the parameter space of the equivalent width (EW), the blue- and red-half-width at half-maximum (HWHM) of the Mg II 2800 Å line, the C IV 1549 Å line, and the C III] 1908 Å blend in samples of Broad Absorption-Line (BAL) and non-BAL quasars at redshift 1.6–2.1. Using this method, we successfully recover correlations well-known in the UV regime such as the anti-correlation between the EW and blueshift of the C IV emission line and the shape of the ionizing Spectra Energy distribution (SED) probed by the strength of He II and the Si III]/C III] ratio. We find this to be particularly evident when the properties of C III] are used to find the clusters, while those of Mg II proved to be less strongly correlated with the properties of the other lines in the spectra such as the width of C IV or the Si III]/C III] ratio. We conclude that unsupervised clustering methods (such as K-means) are powerful methods for finding “natural” binning boundaries in multidimensional datasets and discuss caveats and future work.

Key words: (galaxies:) quasars: emission lines, absorption lines

1 INTRODUCTION

Among the different types of Active Galactic Nuclei (AGN), quasars stand out as the most luminous with typical bolometric luminosities $> 10^{46} \text{ erg s}^{-1}$. Accretion onto a super-massive black hole is now accepted as the main mechanism that powers quasars ([Shakura & Sunyaev 1973](#)). Despite having remarkably similar spectra that typically show a distinct strong blue continuum and broad emission lines with velocity widths of 1000s km s^{-1} , quasar spectra exhibit subtle but strong trends among their emission lines and continua that persist over a wide range of wavelengths and luminosities. These repeated patterns allow us to probe the complex physical processes taking place in quasars’ inner regions.

For example, the profiles of some of the broad high ionization lines (e.g., C IV $\lambda 1550$) exhibit some structure which

indicates that in addition to the Doppler broadening of the line by the central black hole’s gravitational field, a non-virial component of motion in the gas (inflow or outflow) is at work creating this structure seen as blue asymmetries or blueshifts of the peak from the systemic redshift measured in lower ionization lines (such as Mg II, e.g., [Richards et al. 2002](#)). Another notable example is the Baldwin Effect ([Baldwin 1977](#)) seen as an anti-correlation between the strength of C IV $\lambda 1550$ (measured by its equivalent width, EW) and the continuum luminosity at 1550 Å.

More recently it has been demonstrated that C IV EW is anti-correlated with its blueshift and both of those quantities are tied to the X-ray hardness of the quasar in a sense that objects that are soft in the X-ray (characterized by the spectral index α_{ox} ¹) are likely to have weaker C IV and larger blueshifts (e.g., [Richards et al. 2002](#); [Leighly & Moore](#)

¹ The spectral index $\alpha_{ox} = 0.384 \log(f_{2\text{keV}}/f_{2500})$, measures the slope of the flux densities at 2 keV and rest frame wavelength

[★] atammour@uwo.ca

2004; Leighly 2004; Richards et al. 2011). Moreover, an intrinsic fraction of $\sim 20\%$ of quasars of optically selected samples display broad ($\Delta V > 2000 \text{ km s}^{-1}$) UV absorption features; broad absorption line (BAL) quasars (e.g., Weymann et al. 1991; Hewett & Foltz 2003). BAL quasars exhibit relatively weak X-ray emission (e.g., Gallagher et al. 2006) and show blueshifted broad absorption troughs with large velocity offsets indicating an absorbing medium moving outward with velocities up to $25,000 \text{ km s}^{-1}$ (e.g., Weymann et al. 1991). A two-component origin of quasar broad-lines can successfully account for many of their observed properties such as the blueshift-EW anti-correlation of high-ionization lines (e.g., C IV) and its connection with the SED hardness in which broad-lines are a combination of emission from the accretion disk and emission from a fast outflowing gas accelerated by radiation pressure that emerges close to the accretion disk (e.g., Collin-Souffrin et al. 1988; Murray & Chiang 1998; Proga et al. 2000). In this scenario, BAL quasars are seen through the angle covered by the winds and their weaker X-ray emission is a consequence of the winds filtering the energetic photons out of the ionizing radiation (e.g., Proga et al. 2000; Leighly 2004; Leighly & Moore 2004).

1.1 Patterns in Multi-dimensional Space

Looking for systematic patterns appearing repeatedly among the measured variables in large multi-dimensional datasets of quasars can help elucidate the physical driver behind those trends. It is often the case in spectral studies that finding those patterns requires binning the dataset using one or two important parameters then comparing the properties of objects among the different bins by stacking the spectra to create (median or mean) composite spectra (e.g., Croom et al. 2002; Richards et al. 2011; Hill et al. 2014; Shen & Ho 2014; Baskin et al. 2015; Tammour et al. 2015). The binning has been done in most cases in a two-dimensional space with fixed boundaries which follow “traditional” cuts such as 2000 or 4000 km s^{-1} for the width of $H\beta$ (e.g., Sulentic et al. 2007; Tammour et al. 2015) which is not an unreasonable choice as these parameters and boundaries are found to best constrain the properties of the objects under study and to be separating meaningful classes of AGN (e.g., Boroson & Green 1992; Boroson 2002; Sulentic et al. 2000). However, it may be possible to separate objects in a more optimal way when it comes to identifying physical differences.

One notable method used to study multidimensional parameter space in quasar samples is Principal Component Analysis—a statistical method that takes a multidimensional dataset and finds the orthogonal axes (i.e., Eigenvectors) that minimize the variance along each projection (Hastie et al. 2009; Ivezić et al. 2014). This approach was able to uncover interesting correlations among quasar properties (e.g., Boroson & Green 1992; Boroson 2002; Yip et al. 2004). Boroson & Green (1992), for example, found a strong inverse correlation between the strength of the Fe II $\lambda 4570$ complex and the strength of [O III] $\lambda 5007$. This correlation, has been consistently found in other quasar samples and is thought

to originate from the Eddington accretion rate $L_{\text{bol}}/L_{\text{Edd}}$ (e.g., Boroson & Green 1992; Sulentic et al. 2000; Boroson 2002; Shen & Ho 2014).

In this work, we explore the use of unsupervised clustering analysis to find patterns among quasar spectral properties in the UV regime. Unlike supervised learning which uses a labelled training set to assign labels to unlabelled data, unsupervised learning does not require previous knowledge of labels—it is “learning without a teacher” (Hastie et al. 2009, Ch. 14). Clustering analysis is one of the unsupervised learning techniques that aims to find structure in multidimensional data space (Bishop 2009; Hastie et al. 2009). In unsupervised clustering, the data is assigned to clusters according to a chosen metric of similarity such as the Euclidian distance (e.g. the K-means algorithm) or non-Euclidian metrics (e.g., Density Based Spatial Clustering, DBSC) (e.g., Hastie et al. 2009; Ivezić et al. 2014). Despite its robustness to reveal important properties about quasars, PCA assumes orthogonality of the axes defined by the eigenvectors that are a linear combination of the input parameters. This makes interpreting the eigenvectors often physically rather non-intuitive. Clustering of sources by K-means, on the other hand, more closely follows the divisions inherent to the data that we might hope are ultimately due to physics (after accounting for any selection effects).

In §2.1 we describe our samples and in §2.2 we discuss in detail the K-means algorithm and how we apply it to our quasar samples. §3 includes our results and in §4 we discuss some of the caveats of this technique. A summary and conclusion are given in §5.

Throughout this work we use: $H_0 = 70 \text{ km s}^{-1}$, $\Omega_M = 0.3$, $\Omega_\Lambda = 0.7$ (Spergel et al. 2003).

2 ANALYSIS

2.1 Sample Selection

Our starting point in selecting this sample is the Pâris et al. (2014) Quasar Catalog of the 10^{th} Data Release of the Sloan Digital Sky Survey (SDSS-DR10). The large size of this catalog and its improved measurements from those of the SDSS automated pipeline make it ideal for our statistical study of quasar UV spectra. The catalog contains 166,583 quasars with measurements of emission and absorption lines such as the amplitude (AMP), equivalent width (EW), full-width at half-maximum (FWHM) and the red and blue half-width at half-maxima (RHWHM and BHWHM) of Mg II, the C III] blend and C IV. These measurements were obtained with a linear combination of a set of principal components (eigen-spectra) to fit the the full spectrum in a similar fashion to the DR9 Quasar Catalog (Pâris et al. 2012). We choose to restrict the selection to objects within the redshift range of 1.6 to 2.2 (using the PCA redshift values in the Pâris et al. (2014) estimated from fitting for the whole spectrum) to allow us to study some of the interesting UV features such as: Ly α $\lambda 1215$, Si IV $\lambda 1396$, C IV $\lambda 1550$, the He II $\lambda 1640$ and O III] $\lambda 1664$ blend, the Al III $\lambda 1857$, Si III] $\lambda 1892$ and C III] $\lambda 1908$ blend, and Mg II $\lambda 2800$ (Vanden Berk et al. 2001). We note that the redshift values listed in the catalogue could potentially contain inaccurate estimates. This is mainly because high and low ionization lines are commonly

2500 \AA . The quantity 0.384 is the logarithm of the ratio of the frequencies at which the flux densities are measured.

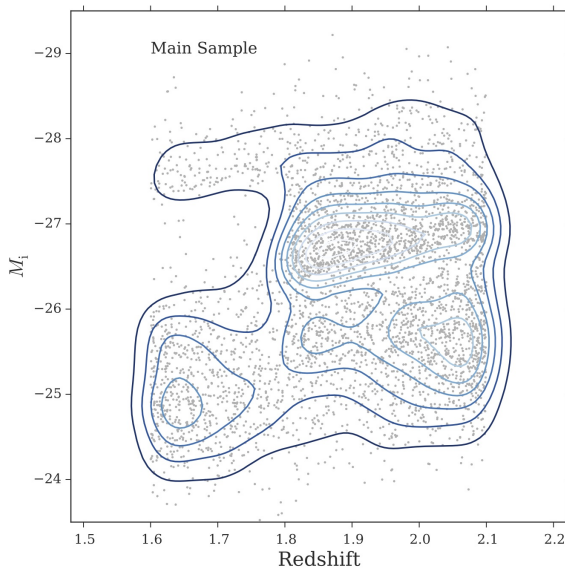


Figure 1. Distribution of M_i vs. z for the objects in our sample. The objects in our main sample are shown with grey dots and blue contours. Given the narrow range of redshift and the SDSS selection, the luminosity distribution is relatively flat with z .

seen to be shifted (from systemic velocities) in different proportions. However, in case the redshifts values are systematically wrong, this will lead to biased measurements that will reveal themselves as separate clusters that can potentially lead to new insights into the physical properties of these objects as we discuss in §3.

We select the initial sample (main sample, §3.1, 3.2, and 3.3) based on the following criteria: EW is > 0 for each of the lines C IV, C III] and Mg II, the EW error is $< 10\%$ of the measured EW for all of the three lines, the BAL flag from visual inspection (BAL_FLAG_VI) is set to 0 (no BAL quasars), and finally we use a signal-to-noise minimum of $S/N > 3$ at 1700 \AA .

This selection gives us a sample of 4110. Fig. 1 shows the distribution of M_i vs. redshift in the main sample. We visually examine the individual spectra of these objects and remove 7 objects with bad spectra (missing flux) and one object that looked like a misidentified BAL quasar. In addition to this main sample, we select two other samples: mixed sample with the EW > 0 and the BAL flag condition relaxed (contains 6463 objects, see §3.4), and BALQ sample with the BAL quasars *only* (contains 1533 objects, see §3.5).

We finally note that there were many cases of heavy narrow absorption in C IV and sometimes in Mg II that were not flagged as BAL quasars in the catalog. We decide to keep those objects in our samples because we wanted to test if the algorithm is able to isolate those objects from the larger quasar population; see §3 for more discussion about this point.

2.2 Unsupervised Clustering with the K-means Algorithm

K-means is one of the most widely used unsupervised clustering algorithms (Hastie et al. 2009; Bishop 2009). Part of the popularity of this algorithm is due to its simple and intuitive design. K-means aims to minimize the “inertia criteria” (i.e., within-cluster sum of squares) given by:

$$J = \sum_{n=1}^N \sum_{k=1}^K \min(||x_n - \mu_k||^2), \quad (1)$$

where μ is the mean in cluster k , N is the total number of data points (samples) and K is the number of clusters. In a nutshell, the algorithm starts by assigning a fixed number of random points to serve as centroids for the clusters to be found. It then proceeds to assign data points that are closer to each centroid to a cluster. Those new cluster members are then used to calculate a new mean which becomes the new cluster centroid. The distances between the new centroids and the cluster members are calculated again and the members are reassigned to the clusters with the smallest distance to their centroids. The algorithm continues to iterate this calculating and reassigning process until it converges. A good visual illustration of this process is given in Bishop (their Fig. 9.1; 2009).

To perform the K-means clustering, we use `scikit-learn`² (Pedregosa et al. 2011) applied to the EW, BWHM, and RWHM measurements provided by the Pâris et al. (2014) catalogue for each of the C IV, C III], and Mg II lines separately. After extensive experimentation, we chose to focus on these three parameters (features) because they adequately describe the lines in terms of strength and structure (asymmetry). In the case of the C III] blend, the blue- and red-HWHM also serve as a measure of the strengths of the Si III] and Al III lines blue-ward of and often heavily blended with the C III] line (we discuss this further in §3.3).

The only parameter that this algorithm requires beforehand is K : the number of clusters that we want (expect) our data to be grouped into. Determining this K is not necessarily a clear-cut exercise as the “ground truth” in our case is not known, i.e., the data points are not labeled. However, heuristic approaches coupled with a good understanding of the dataset and its features are sufficient to determine a range of optimal values of K . We fully expect that the clustering parameters have continuous distributions, perhaps with breaks due more to selection effects than anything. What the clustering does here, is to provide a quantifiable and less arbitrary way of breaking those objects up into groups that can be used to follow the trends in that continuum.

In Fig. 2, we show two metrics that describe the “goodness of fit” for different values of K starting from $K = 2$ to $K = 8$. The top panel gives the sum of the square distances (errors) between data points (samples) and their cluster centroid for each attempt. K at the “elbow” of Fig. 2 (top panel) is the one that minimizes the inertia given in equation 1.

² Using the `sklearn.cluster.KMeans` package: <http://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html#sklearn.cluster.KMeans>

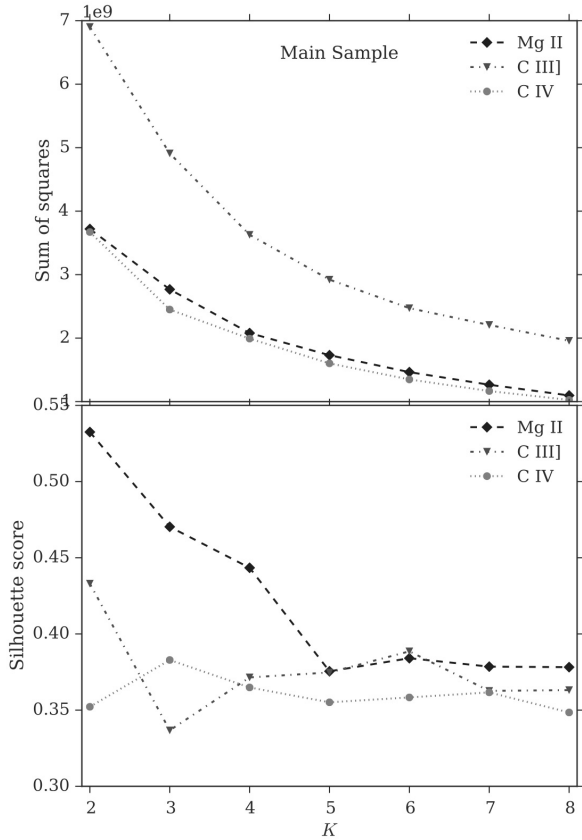


Figure 2. Top: The sum of the squared distances between data points and their cluster centroid (eqn. 1) calculated for $K = 2$ to 8 and using EW, RHWDM, and BHWDM for each of the C III], C IV and Mg II line (blend) separately. For our samples, a range of 3 to 6 is sufficient to capture the curve minimum. Bottom: Average Silhouette scores (eqn. 2) for all data points calculated for the clustering done on each of the three emission lines separately using $K = 2$ –8. Higher scores indicate less overlap among clusters. For Mg II, $K = 2, 3$ and 4 give the best separation. While C III] requires $K = 2$ or $K > 4$ and for C IV $K > 3$ gives the best separated clusters.

Ideally, one can get a good estimate of K by looking for the point where the curve bends (forms an elbow), but in most “real” datasets the curve is smooth and one can only find a range of possible K values that minimizes the square error. The lower panel gives the Silhouette Score, a metric that measures the isolation of clusters by comparing the mean of the distances between a sample and its fellow cluster members, a , and the mean of distances between this sample and the ones in the nearest neighbour cluster, b (Pedregosa et al. 2011; Han et al. 2012):

$$s = \frac{b - a}{\max(a, b)}. \quad (2)$$

This score, s , ranges from 1 for a dense, well-defined cluster to -1 for a less dense cluster that overlaps with the ones surrounding it.

Fig. 2 suggests that 3 or 4 clusters are reasonable for C IV and Mg II while more clusters might be needed for C III]. We proceed to use $K = 3$ to 6 for all 3 lines and we discuss in §3 cases where adding an extra cluster sometimes serves to isolate outliers and create a “cleaner” cluster. Table 1 shows a breakdown of the clustering results for each emission line or blend for $K = 3, 4, 5$ and 6.

Table 1. Number of objects in each cluster formed with K-means applied in the main sample to measurements of EW, BWHM, and RWHM of C III], C IV, and Mg II separately using $K = 3, 4, 5$ and 6 . For each cluster, we also give the coordinates of the cluster centroids in the (EW,BWHM,RWHM) space. The labels “a”, “b”, “c”, etc. refer to the same clusters and composites shown for the rest of this paper. Clusters are ordered according to the increasing values of BWHM of their centroids. Clusters/composites are labeled with the name of the line used for the clustering and the number of clusters used. For example, C III]-5d is cluster d from the clustering done using $K = 5$ in the C III] parameter space. In Appendix B we show that the algorithm is able to find the same unique centroids for the clusters after up to 50 repeats of running K-means.

Line	K	Number of objects EW (\AA), BWHM (km s^{-1}), RWHM (km s^{-1})					
		a	b	c	d	e	f
C III]	3	1437 (24.35,1783.97,2534.25)	2309 (26.85,2754.39,3821.80)	356 (29.68,3376.13,6692.33)	– –	– –	– –
		980 (23.76,1711.87,2226.60)	1919 (26.75,2177.15,3857.84)	328 (29.66,3230.03,6828.22)	875 (26.50,3668.96,3455.78)	– –	– –
	5	392 (22.90,1489.56,1471.68)	1556 (25.44,1988.50,3133.32)	1059 (27.47,2365.86,4354.74)	288 (29.67,3268.35,7016.40)	807 (26.46,3718.25,3425.02)	– –
		375 (22.91,1473.23,1437.42)	1468 (25.40,1937.23,3126.74)	1063 (27.52,2345.50,4337.00)	209 (29.65,2770.46,7281.64)	774 (25.91,3460.89,3191.11)	213 (28.90,4591.70,5247.12)
	6	1713 (43.39,1746.90,1958.73)	1313 (31.04,2044.52,3507.01)	1076 (38.35,3304.51,2698.20)	– –	– –	– –
		1284 (43.97,1562.68,1880.82)	1270 (30.62,1971.40,3453.93)	1020 (41.03,2775.14,2248.06)	528 (36.29,3582.87,3349.97)	– –	– –
C IV	4	977 (45.94,1527.30,1688.75)	1280 (33.91,1945.03,2872.95)	486 (29.37,1999.76,4103.41)	791 (42.30,2892.92,2122.82)	568 (35.79,3504.51,3279.19)	– –
		856 (46.92,1507.34,1617.35)	1060 (34.27,1769.15,2746.57)	417 (28.99,1923.20,4188.69)	756 (33.70,2683.73,3194.66)	730 (43.22,2859.02,2060.24)	283 (37.99,4060.17,3185.96)
	5	2503 (38.14,1795.43,1739.41)	1258 (46.39,3129.40,2591.14)	341 (51.15,3946.39,4451.96)	– –	– –	– –
		2176 (37.42,1757.32,1630.73)	1271 (44.83,2585.56,2628.85)	241 (51.47,3629.85,4850.80)	414 (49.41,4327.26,2591.17)	– –	– –
	6	1468 (35.36,1589.28,1498.43)	1371 (42.53,2262.72,2087.34)	707 (46.65,2818.55,3085.96)	159 (53.23,3799.77,5284.70)	397 (49.38,4393.77,2639.92)	– –
		1482 (35.42,1596.65,1496.80)	1366 (42.54,2235.09,2131.36)	622 (46.98,2830.91,3192.47)	77 (51.61,3100.29,6065.12)	366 (48.02,3989.87,2222.78)	189 (52.35,4747.26,3891.82)
Mg II	3	2503 (38.14,1795.43,1739.41)	1258 (46.39,3129.40,2591.14)	341 (51.15,3946.39,4451.96)	– –	– –	– –
		2176 (37.42,1757.32,1630.73)	1271 (44.83,2585.56,2628.85)	241 (51.47,3629.85,4850.80)	414 (49.41,4327.26,2591.17)	– –	– –
	5	1468 (35.36,1589.28,1498.43)	1371 (42.53,2262.72,2087.34)	707 (46.65,2818.55,3085.96)	159 (53.23,3799.77,5284.70)	397 (49.38,4393.77,2639.92)	– –
		1482 (35.42,1596.65,1496.80)	1366 (42.54,2235.09,2131.36)	622 (46.98,2830.91,3192.47)	77 (51.61,3100.29,6065.12)	366 (48.02,3989.87,2222.78)	189 (52.35,4747.26,3891.82)
	6	2503 (38.14,1795.43,1739.41)	1258 (46.39,3129.40,2591.14)	341 (51.15,3946.39,4451.96)	– –	– –	– –
		2176 (37.42,1757.32,1630.73)	1271 (44.83,2585.56,2628.85)	241 (51.47,3629.85,4850.80)	414 (49.41,4327.26,2591.17)	– –	– –

2.3 Median Composite Spectra

To better visualize the results of the clustering and to examine the properties of the other lines which were not used in the clustering analysis, we create median composite spectra from the objects in each of the clusters we find. We start by correcting all individual spectra for Galactic extinction and redshift using the extinction at the g-band and the redshift values estimated from PCA as quoted in the [Pâris et al. \(2014\)](#) catalogue. We normalize the individual spectra using the median of the flux between 2360 Å and 2390 Å and apply 3σ clipping³ before they are median-combined.

3 RESULTS AND DISCUSSION

3.1 Mg II Clusters

The Mg II clusters generated using EW, BWHM and RWHM of Mg II are shown in Fig. 3 for $K = 3, 4, 5$, and 6. The figure shows the 2-D plane of the blue- and red-HWHM as the x and y -axes. The figure also shows the mean EWs for each cluster. While the EW is still contributing to the clustering (see values of the cluster centroid coordinates in Table. 1), the two width parameters have more contribution in determining the clustering results. In Fig. 4, we show the $K = 5$ run only along with the BALQ sample quasars (not used in the clustering analysis) over-plotted and the distributions of both the full main and BALQ samples. The Mg II line is largely symmetric (with the exception of Mg II-e5 which we discuss below) and its average EW increases in the clusters with its width as Fig. 4 shows. The red- and blue-HWHM of Mg II appear to have similar distributions in both the BAL and non-BAL quasar samples.

We median-stack objects in each one of the clusters shown in Fig. 3 and show the results in Fig. 5. The projected width of Mg II is the predominant quantity driving differences among clusters (as the Mg II profiles indicate) and is potentially probing differences in the black hole masses of objects in each cluster. However, with this change in the widths of Mg II, we do not see any clear corresponding changes in the strengths, widths or asymmetries of any of the other lines in the composites. C IV for example appears to have nearly identical profiles in all the composites regardless of the width of Mg II. Because of the small range of luminosity in the sample, the Mg II clusters are likely reflecting BH masses. This lack of correlation between the line widths of Mg II and C IV is consistent with the known discrepancy in the estimates of black hole masses frequently done using single epoch C IV measurements (e.g., [Baskin & Laor 2005](#); [Shen & Liu 2012](#)).

The absorbed peak of C IV in composite Mg II-d5 (Fig. 3) appears to be a result of many narrow absorption lines in C IV in this cluster that are evident in the median composite due to the low number of objects (159 objects). We visually examine the spectra of individual objects in this cluster and find that most of them have very broad Mg II profiles and frequent narrow absorption in Mg II and/or

C IV. Most objects in cluster Mg II-e5 in Fig. 4 fall below the diagonal causing the skewed blue-ward Mg II profile in Fig. 5. Visual examination of objects in this cluster showed a few cases where the redshift might have been poorly estimated potentially due to narrow absorption in C IV and/or Mg II which likely causes the automated fitting to misidentify the line peak. Most objects, however, appear to have C IV and C III] peaking at systemic redshift. We check for discrepancies in the redshift estimates by comparing the PCA redshift (which we used in shifting our spectra to rest-frame) to those estimated from Mg II alone, as given in the [Pâris et al. \(2014\)](#) catalog, and found that they are in good agreement (Table 2). Indeed Mg II is known to be one of the most reliable redshift indicator among broad permitted lines in quasar spectra (e.g., [Hewett & Wild 2010](#)). Our visual inspection also identified an enhancement red-ward of Mg II in many of the objects in the Mg II-e5 cluster. Such enhancement can be a result of stronger FeII features blended with Mg II which could cause the observed blue-skewness of the Mg II line.

3.2 C IV Clusters

We repeat the same analysis done in (§3.1) but this time using the C IV parameters as the features fed to K-means with $K = 3$ and up to 6. In Fig. 6, we show the clusters for the $K = 5$ case with the BAL quasars over-plotted and the distributions of their blue- and red-HWHM compared to the distributions of the main sample. The C IV clusters are significantly more compact than those of Mg II with a dynamical range of ~ 1000 – 5000 km s^{-1} for the blue- and red-HWHM and the average C IV EWs are larger in objects with narrower C IV (Table 1).

The blue-HWHM of the main and BALQ samples appear to follow similar distributions (top panel in Fig. 6) but the distribution of the red-HWHM in the BALQ sample (right panel of Fig. 6) is skewed towards higher values than that of the main sample, i.e., the red-HWHM of C IV is larger in most BAL quasars than that of non-BAL quasars. We point out, however, that the C IV emission-line measurements for BAL quasars can be sometimes unreliable because the BAL features often interfere with the emission-line fits as the absorption can in many cases be fully or partially superimposed on the emission-line which undermines the measurements of the emission line.

Notably, Fig. 7 shows that despite using C IV properties *only* to generate the clusters, the other emission lines in the composite spectra generated from the clusters are corresponding to C IV in a way that is expected from observations of quasar UV spectra. For example, composite C IV-a5 in Fig. 7 has narrow, symmetric C IV which shows no shift from the systemic redshift and has a smaller Si III]/C III] ratio and a prominent He II (also strong Si IV, OI and SiII and narrower Mg II), while the composites with blueshifted and weaker C IV (e.g., composite C IV-c5) have larger Si III]/C III] ratios and weaker He II. We point out here that the presence of a strong He II line coupled with a higher Si III]/C III] ratio is an indication of a harder ionizing SED –although other interpretations of a high Si III]/C III] ratio alone (such as a higher density) are also possible ([Leighly & Moore 2004](#); [Casebeer et al. 2006](#)). These observations are in line with the disk-wind model of the broad-lines in quasars

³ Each median spectrum only includes points contributing to the spectrum that are inside the range of 3σ standard deviations in each wavelength bin.

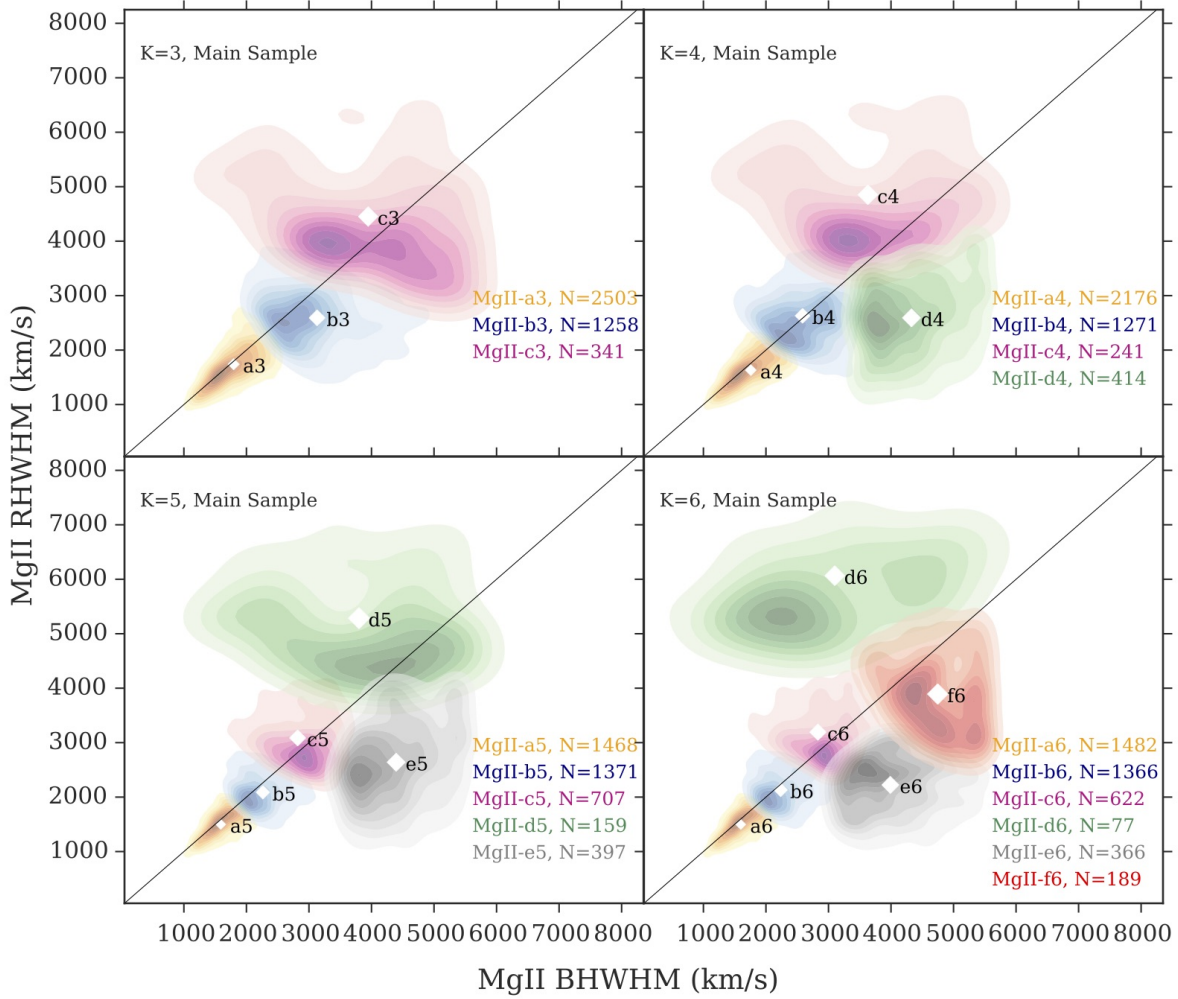


Figure 3. Red- vs. Blue-HWHM plane for clusters found using $K = 3$ to 6 (as labeled in the top-left corner of each panel) and the EW, BHHWM and RHHWM measurements of Mg II. The cluster labels match the ones given in Table 1. The white diamonds show the projected locations of the cluster centroids in this plane and their sizes are proportional to the mean EW of each cluster. The diagonal black lines mark the 1 to 1 ratio (symmetric line).

which predicts that blueshifted high ionization lines are a result of a fast-moving wind emerging from the accretion disk and accelerated outward by radiation pressure, while intermediate and lower ionization lines are emitted closer to the base of the wind at the accretion disk and are therefore seen to be symmetric and at systemic redshift (Leighly 2004; Leighly & Moore 2004; Richards et al. 2011). The shape of the ionizing SED in the disk-wind model plays a major role in the contributions of each of the disk and wind components to the broad emission lines; quasars with higher UV luminosity relative to the X-ray have a stronger wind component (more blueshift in C IV), while quasars with harder X-ray over-ionizes the gas and have therefore less contribution from the wind to the broad lines (e.g., Casebeer et al. 2006; Leighly & Moore 2004; Leighly et al. 2007; Kruczek et al. 2011; Richards et al. 2011).

We also see in Fig. 7 that clustering done using C IV properties results in Mg II profiles that are rather similar regardless of the C IV profiles (similar to what we see in the Mg II clusters in Fig. 5). This again cautions that estimates of black hole masses using C IV could result in values that do not necessarily reflect the ones estimated by the less biased Mg II. We finally point out that the weak absorption seen in the C IV profiles of composites C IV-d5 and C IV-e5 might indicate a higher number of narrow absorption among objects in these two clusters. This feature is unlikely to be due to low numbers as cluster C IV-c5 has fewer objects but its C IV profile shows no sign of absorption. Fig. 6 shows that these two clusters include objects with extremely high BHHWM. We visually examine the spectra of objects from clusters C IV-d5 and C IV-e5 with the highest BHHWM and find that most of these objects have either absorption

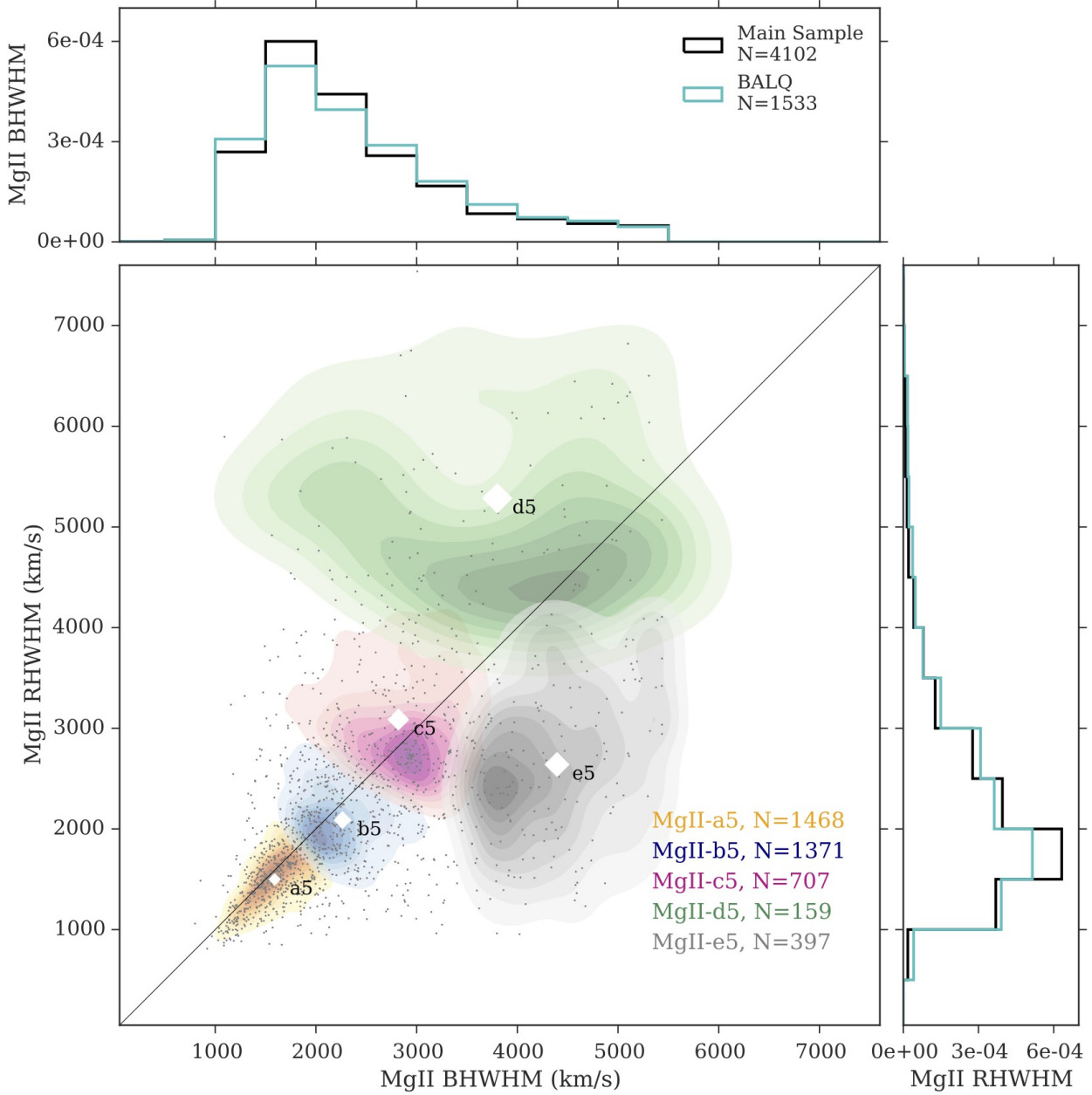


Figure 4. Red- vs. Blue-HWHM of Mg II for the $K = 5$ clustering run. We over-plot the BAL quasars (not included in the main sample) as grey dots and show the fractional distributions of the Mg II blue- and red-HWHM of the main and BALQ samples in the upper and left panels respectively. It is clear that Mg II has a fairly symmetric profile except for the d5 and e5 clusters which we discuss further in the text. Both the blue- and red-HWHM of Mg II of the BAL quasars have distributions similar to those in the main sample. The white diamonds show the projected locations of the cluster centroids and their sizes correspond to the mean EW for each cluster. The average EW of Mg II increases with its width as the sizes of the white diamonds appear to increase diagonally with the RWHHM and BWHHM.

in their C IV profiles that is not flagged as BAL or have C IV profiles that are highly skewed blue-ward (could be a blueshift or absorbed flux on the red side). Table 2 indeed shows that the C IV redshifts in these two clusters slightly diverge from the PCA redshifts. In Fig. 8 we show examples of C IV profiles for four objects with extremely high BWHM in clusters C IV-d5 and C IV-e5. We see no clear

trend in these two clusters with respect to the properties their SED as seen by the He II and the Si III]/C III] ratio in the corresponding composites.

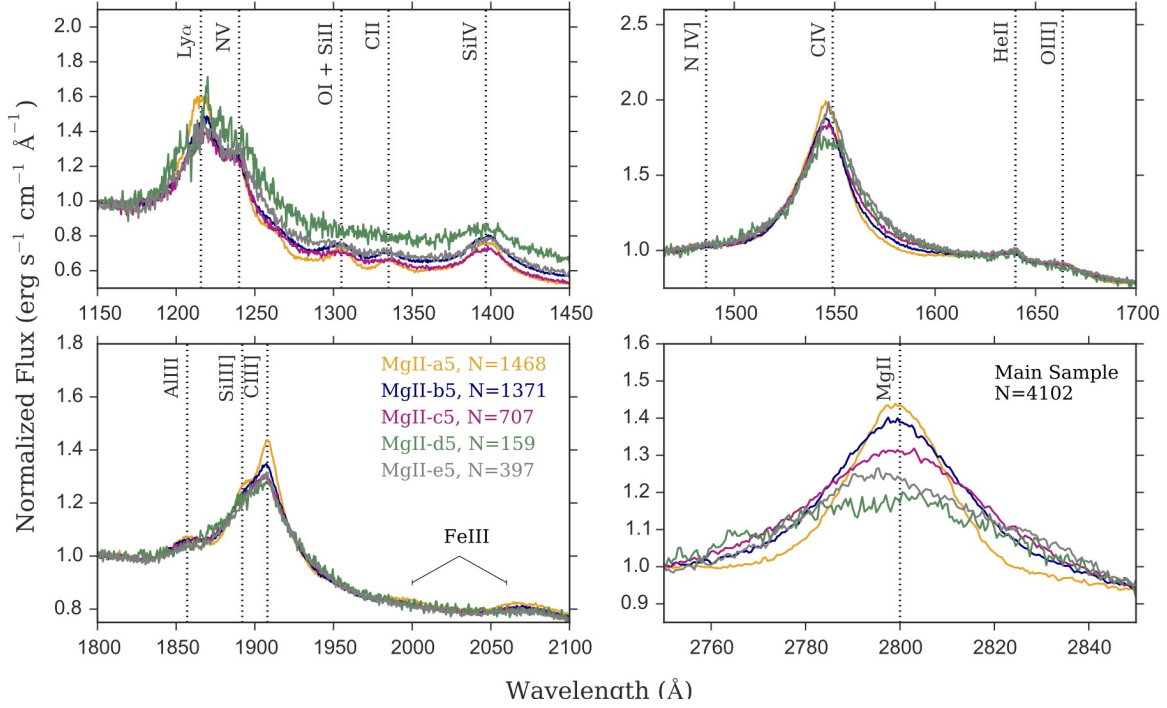


Figure 5. Median composite spectra made from the objects in the Mg II clusters with $K = 5$ shown in Fig. 4. The profiles have been normalized locally at the starting wavelength of each panel. The numbers in the lower-left panel refer to the number of objects in each composite/cluster (see also Table 1). The width of Mg II appears to be the main driver for the clustering and is potentially probing average black hole masses in each cluster. Composites Mg II-d5 and Mg II-e5 have objects with relatively high red- and blue-HWHM respectively as shown in Fig. 4. The weak dip in the C IV profile in composite Mg II-d5 betrays the presence of narrow absorption features in many objects in this cluster.

Table 2. Comparison between the redshifts estimated from PCA and individual lines as given in the Pâris et al. (2014) catalog for each of the clusters we identify in the Mg II, C IV and C III] runs. We report the statistic and p value of the K-S test which tests the null hypothesis that the two samples are drawn from the same distribution. Redshift estimates using Mg II have similar distributions to those from PCA with high significance in all clusters. Redshifts from C III] are slightly less identical to those from PCA except for C III]-e5 which appears to have significantly different redshifts from those estimated from PCA. C IV redshifts vary in their similarities to those from PCA but are mostly in good agreement with them.

Cluster	Num Obj	$\Delta Z(\text{Mg II})$	p	$\Delta Z(\text{C III]})$	p	$\Delta Z(\text{C IV})$	p
C III]-a5	392	0.020	1.000	0.018	1.000	0.015	1.000
C III]-b5	1556	0.008	1.000	0.013	0.999	0.019	0.932
C III]-c5	1059	0.009	1.000	0.021	0.975	0.029	0.749
C III]-d5	288	0.028	1.000	0.042	0.960	0.059	0.685
C III]-e5	807	0.020	0.997	0.061	0.098	0.038	0.583
C IV-a5	977	0.011	1.000	0.012	1.000	0.013	1.000
C IV-b5	1280	0.011	1.000	0.023	0.894	0.026	0.784
C IV-c5	486	0.019	1.000	0.035	0.923	0.033	0.952
C IV-d5	791	0.011	1.000	0.028	0.916	0.037	0.655
C IV-e5	568	0.016	1.000	0.046	0.582	0.048	0.533
Mg II-a5	1468	0.012	1.000	0.015	0.996	0.026	0.704
Mg II-b5	1371	0.009	1.000	0.030	0.566	0.028	0.662
Mg II-c5	707	0.018	1.000	0.031	0.879	0.024	0.986
Mg II-d5	159	0.050	0.986	0.038	1.000	0.057	0.956
Mg II-e5	397	0.038	0.935	0.033	0.982	0.028	0.998

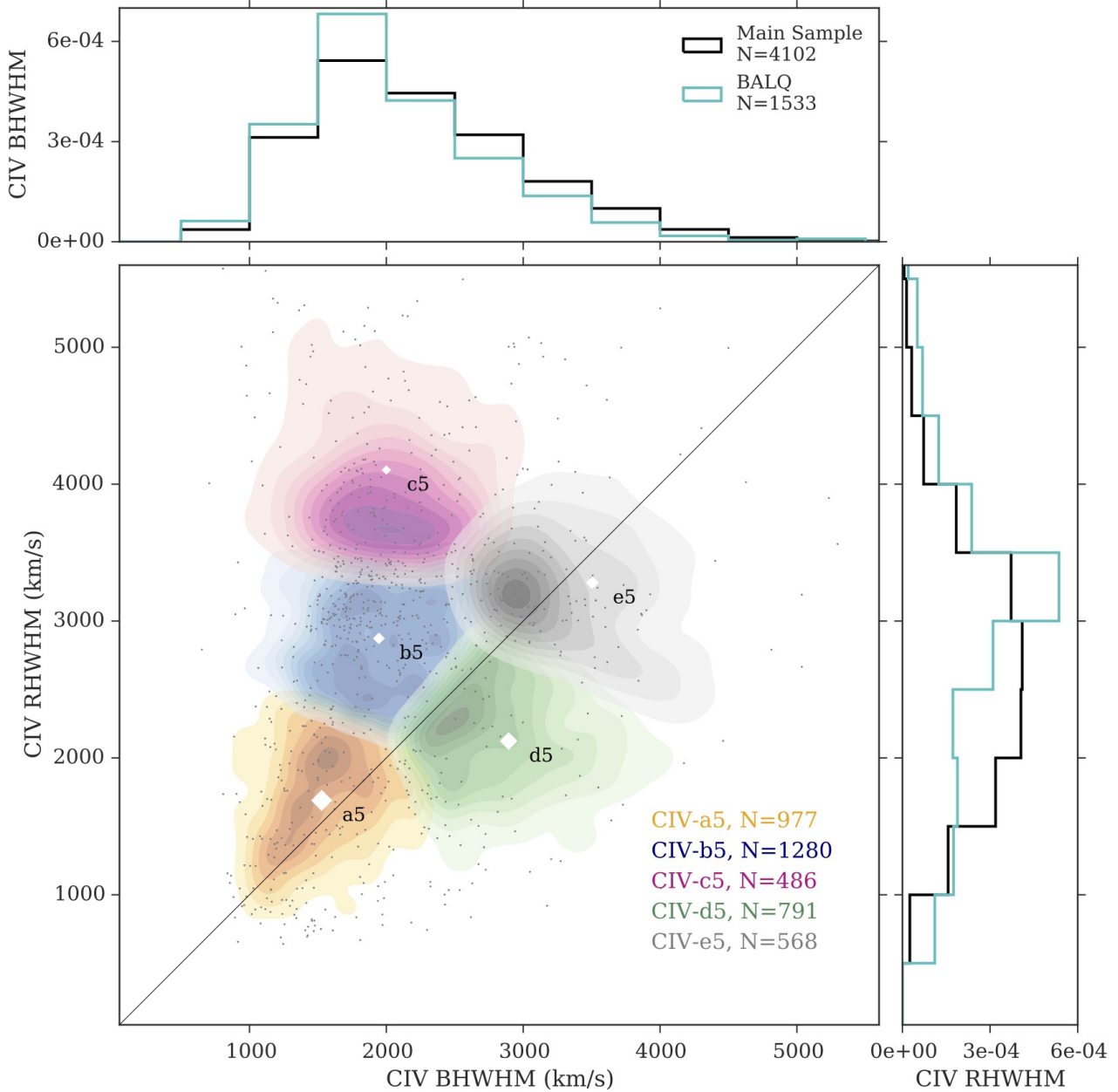


Figure 6. Red- vs. Blue-HWHM of C IV for the $K = 5$ clustering. We overplot the BAL quasars as grey dots and show the fractional distributions of the C IV blue- and red-HWHM of the BALQs and the full sample in the upper and left panels respectively. The distribution of the red-HWHM in the BALQ sample (right panel) peaks at higher values than that of the main sample indicating the red-HWHM of C IV is larger in most BAL quasars compared to that of non-BAL quasars. We note that the BALQ sample contains many objects with measured HWHM of -1. We discuss these objects further in §3.4 and 3.5 and show that the algorithm was able to isolate them from the rest of the objects in the sample.

3.3 C III] Clusters

We repeat the clustering analysis here in a similar fashion to the previous two lines and with $K = 3, 4, 5$, and 6 for the EW, BWHM and RWHM of C III]. Fig. 9 shows the results of those runs in the BWHM–RWHM plane for the $K = 5$ case. The figure also shows the distributions of

the blue- and red-HWHM of C III] for the main and BALQ samples and the average EW for each cluster. The blue- and red-HWHM of C III] in BAL and non-BAL quasars are fairly consistent and the EWs of objects in the different clusters do not appear to have a major role in driving the clustering. Furthermore, the C III] measurements in the [Pâris et al. \(2014\)](#) catalog refer to the C III], Si III], and Al III com-

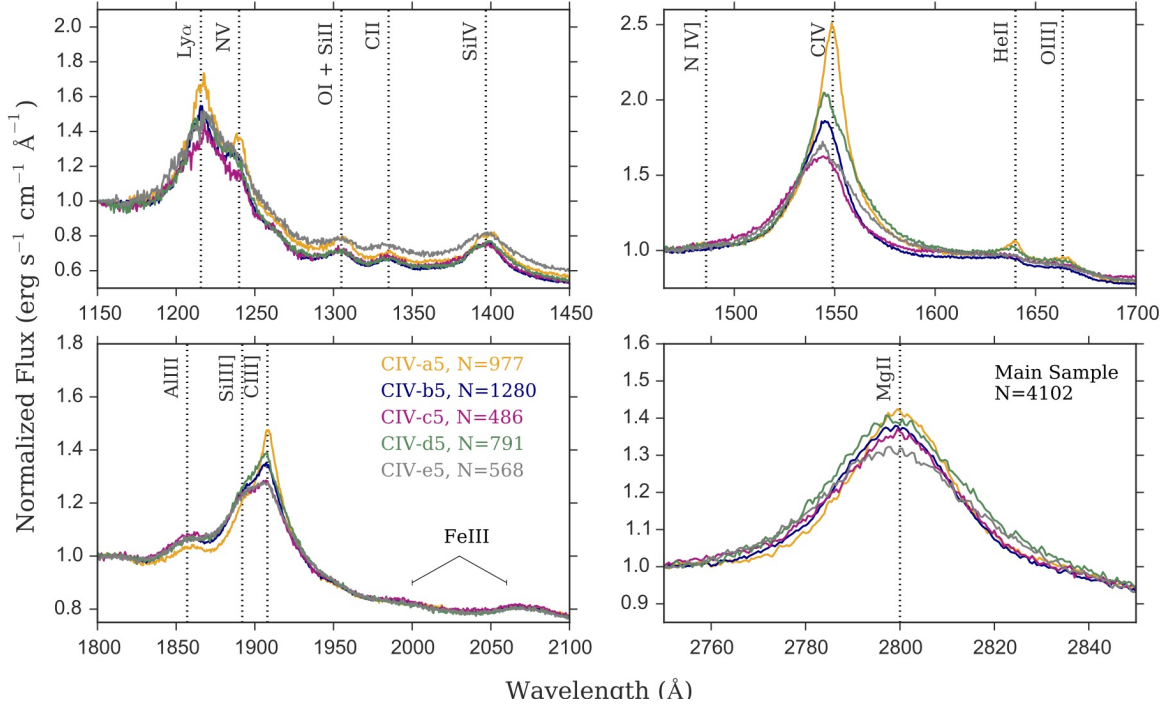


Figure 7. Median composite spectra made from the main sample objects in the C IV clusters with $K = 5$. The profiles have been normalized locally at the starting wavelength of each panel. The labels in the lower-left panel refer to the number of objects in each composite/cluster (see also Table 1). C IV profiles shift from a large EW, strongly peaked and symmetric line (C IV-a5) to a broader, weaker and blueshifted one (C IV-c5) at the extremes. C IV-a5 also shows narrow, peaked profiles in Ly α and C III] with a low Si III]/C III] ratio while C IV-c5 has a higher Si III]/C III] ratio. The prominence of He II and the low Si III]/C III] ratio in C IV-a5 are signs of a hard ionizing continuum relative to C IV-c5. Mg II lines have nearly similar widths regardless of the shape of the C IV. The weak absorption seen in the C IV profiles on C IV-d5 and C IV-e5 is discussed in the text (§3.2 and Fig. 8).

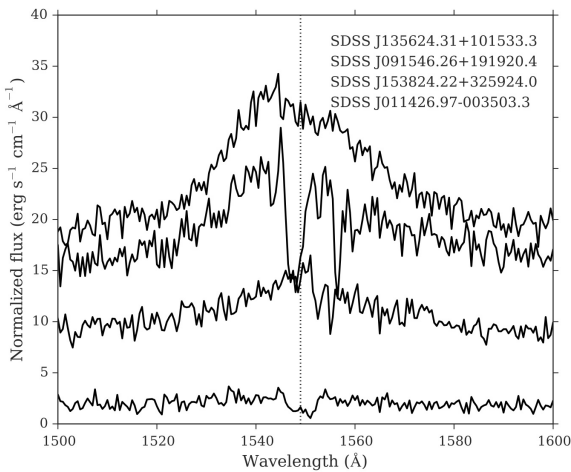


Figure 8. Example spectra of individual objects from the C IV-d5 and C IV-e5 composites in the main sample (Fig. 7). Objects in these clusters have either highly blueshifted C IV and/or have absorption that is missed from the BAL visual inspection or is not broad enough to be classified a BAL quasar.

plex measured without any deblending. This means that the EW of this blend can be similar in two objects but with significantly different contributions from the lines in the blend. The red- and blue-HWHM in the C III] blend in this catalog are measured from the line centroid set at the location of the line peak (maximum intensity) of the blend (Pâris et al. 2012, 2014). This means that objects with extremely high RHHM in Fig. 9 are simply objects with weaker C III] relative to the other two lines in the blend (e.g., cluster C III]-d5 vs. cluster C III]-b5). This becomes clear when we examine the composite spectra in Fig. 10 generated from median-combining objects in each cluster in Fig. 9. A great advantage of using K-means is the ability to use the C III] blend as a whole (the entire line blend) to isolate objects with varying strengths of C III], Si III], and Al III] without having to deblend them (which might indeed be challenging with survey quality data).

Again in this case, the clusters made using the C III] complex properties are reflecting similar trends to those found in the composites in Fig. 7 and are perhaps even more pronounced. Indeed objects with large EWs and narrow, symmetric profiles for C IV (also Ly α , He II, and Mg II; composite C III]-a5 in Fig. 10) have a lower Si III]/C III] ratio while objects with lower EWs and broader blueshifted lines (composites C III]-d5 and C III]-e5 in the same figure) have a higher Si III]/C III] ratio. This suite of observations

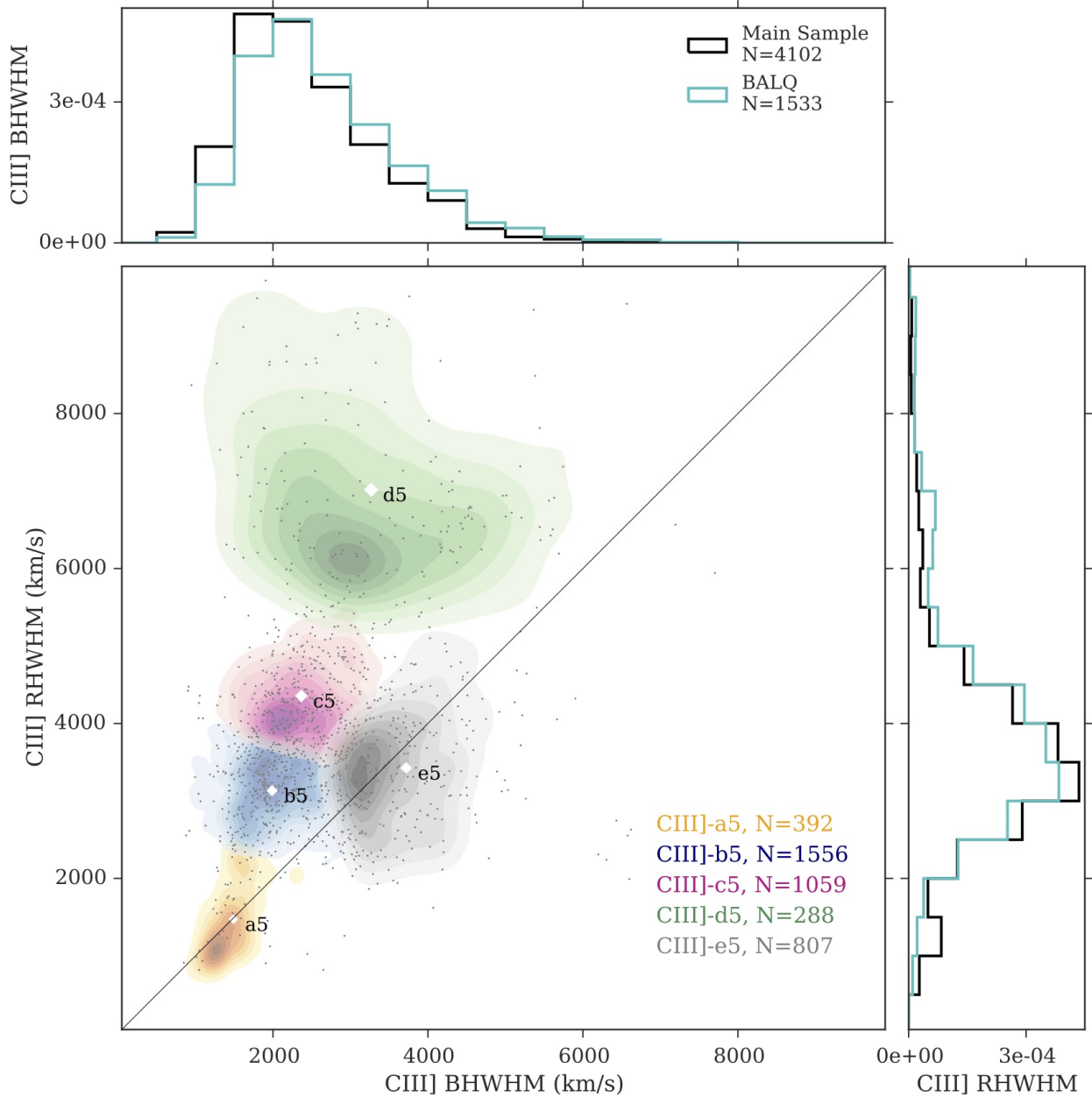


Figure 9. Red- vs. blue-HWHM of C III] for the $K = 5$ clustering in the main sample. The diagonal black line marks to 1:1 ratio. The white diamonds mark the projected locations of the cluster centroids and their sizes correspond to the average EW of the C III] blend of objects in each cluster which appear to be consistent among clusters. We over-plot the BAL quasars (not used in the clustering) with grey dots and show the fractional distributions of the C III] blend blue- and red-HWHM of the BAL quasars in the upper and right panels respectively along with those of the main sample. We see no significant differences in distributions of both samples. The larger values of the red-HWHM of the C III] blend (cluster C III]-d5) reflect a lower C III] line peak relative to the other two lines in the blend as becomes clear when we look at the composites (Fig. 10).

is in agreement with the disk-wind model discuss in §3.2 in which the shape of the ionizing SED has strong effects on the lines in a sense that a quasar with a soft ionizing SEDs (weak He II, high Si III]/C III] ratio, for example C III]-d in Fig. 10) exhibits stronger wind component (weaker C IV with larger blueshifts), while a harder SED over-ionizes the

wind resulting in a stronger C IV that is centred at zero velocity (composite C III]-a; e.g., Casebeer et al. 2006; Leighly et al. 2007; Richards et al. 2011).

Clusters C III]-a5 through C III]-d5 in Fig. 9 appear to follow a sequence with increasing red- and blue-HWHM – a trend that is also reflected in their C III] profiles in Fig. 10.

Cluster C III]-e5 seems however to fall off this sequence with its red-HWHM similar to those of C III]-b5 and C III]-c5 but with larger blue-HWHM. Its composite spectra in Fig. 10 shows that the line profile also falls off the sequence of the decreasing C III]/Si III] ratio from C III]-a5 through C III]-d5. In addition, the C IV profile in this cluster does not follow the gradual increase in blueshift from C IV-a5 to C IV-d5 and its Mg II appears to be slightly skewed blueward systemic. Table 2 shows that the C III] redshift estimates for this cluster are significantly different from those of PCA (p value $\simeq 0.1$) which potentially indicates discrepancy in the blue- and red-HWHM of the C III] blend in this cluster due to problematic redshift estimates. Table 2 also shows that C IV redshifts for this cluster are in less agreement with the PCA redshifts.

Another notable feature in Fig. 10 is the relatively strong N IV] 1486 Å line which becomes prominent in composite C III]-a5. The low Si III]/C III] ratio and the strong He II in this composite indicate a harder ionizing SED. This nitrogen feature is rather uncommon in quasar spectra – Bentz et al. (2004) found less than 1% of their SDSS sample has enhanced nitrogen and our visual examination of selected objects from cluster C III]-a5 confirms the scarcity of N IV] 1486 Å in this cluster. Jiang et al. (2008) studied a sample of 293 quasars with strong nitrogen emission lines (N IV] λ 1486 or N III] λ 1750) and found that quasars with higher nitrogen abundance share similar properties with the overall quasar population though higher nitrogen abundance is associated with higher fraction of radio-loud objects. We look at the fraction of objects detected by FIRST and find that indeed, cluster C III]-a5 includes a slightly higher fraction of quasar detected by FIRST (12.8 %) compared to the other clusters which contain a range of 7.1% to 8.7% FIRST-detected objects.

3.4 Clustering on the Mixed Sample

Next, we set the BAL flag to 1 in the selection process mentioned in §2.1, allowing BAL quasars to be included in the sample. We also remove the constraint that $EW > 0$ and keep the redshift limit (1.6–2.1) and the S/N and EW uncertainty cutoffs: $S/N(\lambda 1700) > 3$, $EW \text{ error} < 10\%$ of measured EW. Adding BAL quasars to the sample increased its size to 6463 (not including the 15 objects we previously removed from the main sample). The purpose of defining this mixed sample is to examine the C III] properties of the entire quasar population (BAL and non-BAL quasars). We do not use C IV and Mg II measurements for the clustering analysis here because (by definition) at least one of these lines shows absorption in BAL quasars and can therefore bias the results.

Fig. 11 shows the clusters using the C III] blend properties in the mixed sample for $K = 6$ run. We find that in the mixed sample, the algorithm is separating a group of objects that were not fit in the Pâris et al. (2014) catalog and a “-1” was entered instead of the measurements of the red- and blue-HWHM and EW. It is only when $K = 6$ that cluster C III]-a6 is almost purely made of the “negative” objects (971 objects and only one object with non-negative values (22.6 Å, 741 km s⁻¹, 695.5 km s⁻¹)). Of these objects in cluster C III]-a6 $\sim 33\%$ are flagged as BAL quasars in the Pâris et al. (2014) catalog while the rest are not marked as

BAL quasars. Visual inspection of the spectra of individual objects in cluster C III]-a6 shows that many of them have “flat” spectra (no emission lines and in some cases what looks like weakly absorbed lines) and so the automated fitting routine was not able to find lines to fit either because the lines are very weak or because of the strong BAL troughs. Objects with BAL troughs could be easier to misfit as the absorption troughs could be highly blueshifted and cause the fitting routine to fail. Some of the objects in cluster C III]-a6 with the weak emission lines might potentially meet the definition of the so-called weak-lined quasars (e.g., Luo et al. 2015; Shemmer & Lieber 2015, and references therein). This is supported by the presence of a more prominent Fe III feature blue-ward of C III] known to be stronger in weak-lined quasars as can be seen in Fig. 12 which shows the composite spectra generated from objects in each of the clusters.

It is also intriguing to see that the objects collected in cluster C III]-a6 (with the -1 values) are generating a weak, strongly blueshifted C IV profile which may contain absorption features as can be inferred from the dip in Ly α . This weak C IV line is also associated with the highest Si III] to C III] ratio indicating a soft ionizing SED on average in those objects compared to the other clusters. Similar to what we see in the main sample (§3.2 and §3.3), the cluster with the narrow, strong C IV line which is centered at the systemic redshift (composite C III]-b6) has the strongest He II and the smallest Si III] to C III] ratio consistent with a hard ionizing continuum which inhibits wind formation.

Finally, it is worth mentioning that the fraction of BAL quasars in each cluster appears to support the notion that BAL quasars are more likely to have softer (X-ray weaker) SEDs as we find that the numbers of BAL quasars in each of the mixed cluster decreases gradually with the Si III]/C III] ratio; cluster C III]-b6 in Fig. 12 (with the strongest He II and lowest Si III]/C III] ratio) has only 11% BAL quasars while clusters C III]-a6 and C III]-e6 (with relatively weaker He II and higher Si III]/C III] ratio) have each 33% of its objects flagged as BAL quasars in the catalog. The percentage in the rest of the clusters is: C III]-c6 has 19% BAL quasars, C III]-d6 has 24% BAL quasars and C III]-f6 has 23% BAL quasars.

3.5 Sample with BAL Quasars Only

We now look at the properties of 1533 BAL quasars selected as described in §3.4 but with the BAL flag set to 1. Here, too, the algorithm is separating the objects with no fitting in the catalog (with -1 entries for the EW and widths). Up to $K = 6$, the cluster includes 8 objects with > -1 values with the remaining 327 objects with -1 for the measurements of C III] EW and blue- and red-HWHM.

Figure 13 shows the median composite spectra generated from objects in the $K = 5$ run. We find that the C IV profiles (both in emission and absorption) are changing with the properties of the C III] blend used to generate the clusters. Composite C III]-b5, for example, with the lowest Si III]/C III] ratio, has a strong, narrow C IV emission-line centred at the systemic redshift and a deep absorption trough with low V_{min} . On the other hand, composites with larger Si III]/C III] ratios (C III]-d5 and e5) have blueshifted C IV emission-lines and higher V_{min} for their absorption troughs. Composite C III]-5a – mostly com-

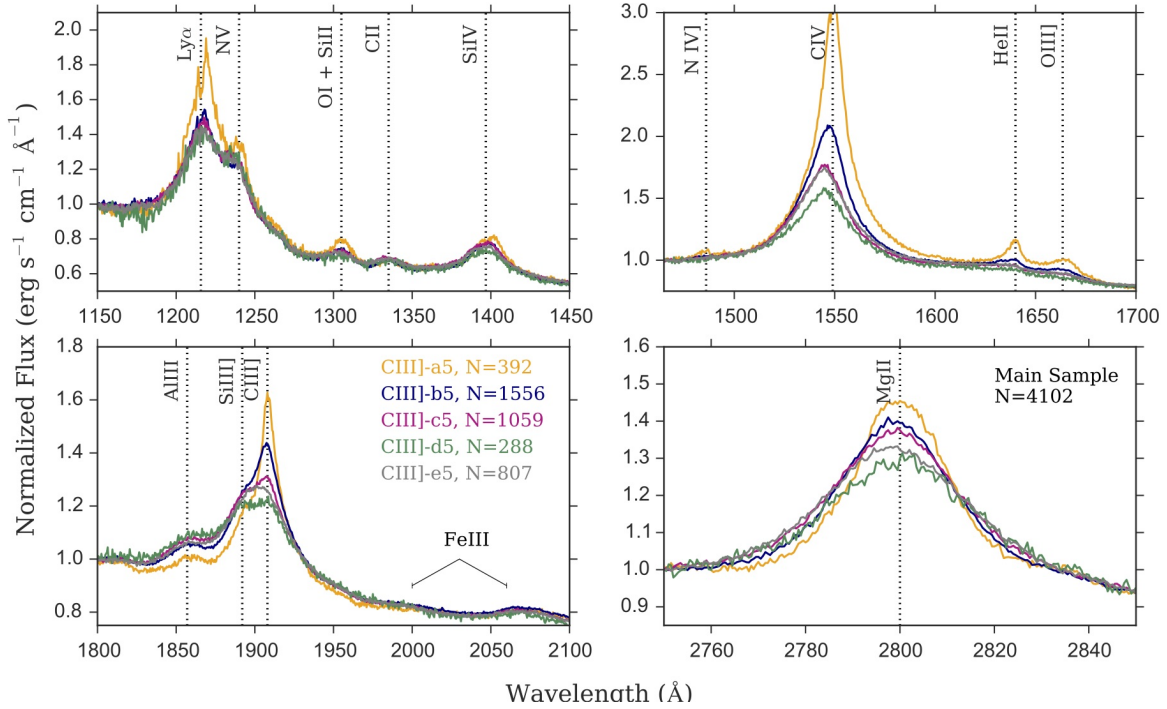


Figure 10. Median composite spectra made from the objects in the C III] clusters with $K = 5$ shown in Fig. 9. The profiles have been normalized locally at the starting wavelength of each panel. The numbers in the lower-left panel refer to the number of objects in each composite/cluster (see also Table 1). Composite C III]-a5 with large EWs and narrow, symmetric profiles for C IV, Ly α , He II, and Mg II has lower Si III]/C III] ratio while composite C III]-d5 has smaller EW and broader blueshifted C IV with higher Si III]/C III] ratio indicating a softer ionizing continuum and a larger wind contribution to the C IV profile.

prised of objects with the -1 measurements of C III]– has the highest Si III]/C III] ratio and shows weak, blueshifted C IV emission-line and a shallower trough and a very high V_{min} . The absorption troughs in Si IV appear to follow similar trends of those seen in C IV. The lack of any absorption showing in the Mg II profiles in those composites indicates that objects in these clusters contain either only high ionization BAL quasars or a mixture of both high and low ionization BAL quasars that averages out in the Mg II profiles.

In future work (Tammour et al., in prep), we explore the properties of the C III] complex and the absorption lines in detail for a sample of BAL quasars.

4 CAVEATS

The K-means algorithm is an unsupervised clustering method. This means that, by definition, there are no fixed labels to the objects in the clusters defined by K-means and it is the investigator’s judgement whether or not to accept the clusters found by the algorithm. The two main decisions to be made beforehand are what features to use to define the parameter space, and how many clusters to group the objects into. Deciding which features to use is limited to what is available and is strongly related to the problem in hand. In this work, we use the EW and the blue- and red-HWHM of individual lines or line blends. We decide to use these

three parameters as they adequately represent the relative strength and structure (asymmetry) of the emission lines. Determining K (the number of clusters), on the other hand, is a heuristic exercise. It requires familiarity with the dataset and the scientific question in hand, and a few iterations of choosing a K value and examining the output. In our case, we find that even though $K = 3$ gives nicely separated clusters, increasing K to 5 or 6 allows us to isolate some of the outliers that are otherwise “hiding” in the clusters. The two metrics we use (discussed in §2.2) are heuristic and are useful for guidance –this is why we opt to test a range of K values rather than one single value. In choosing $K = 5$ (and $K = 6$ for the mixed sample) we aim to balance using a low number of clusters that could potentially hide interesting properties of the objects and fracturing the cluster into smaller bits that might not show interesting properties.

Finally, we note that results based on machine learning techniques strongly rely on the input data fed to the algorithm. For this work we use measurements from the Quasar properties catalog of Pâris et al. (2014). This catalog includes results of automated measurements of 166,583 quasars and cases where the fitting routine failed to find a proper fit or where broad absorption was not properly flagged are not surprising. One particularly useful outcome of using K-means is its ability to isolate objects with special features from a large sample without prior knowledge of their peculiar properties. An example of this is isolating a

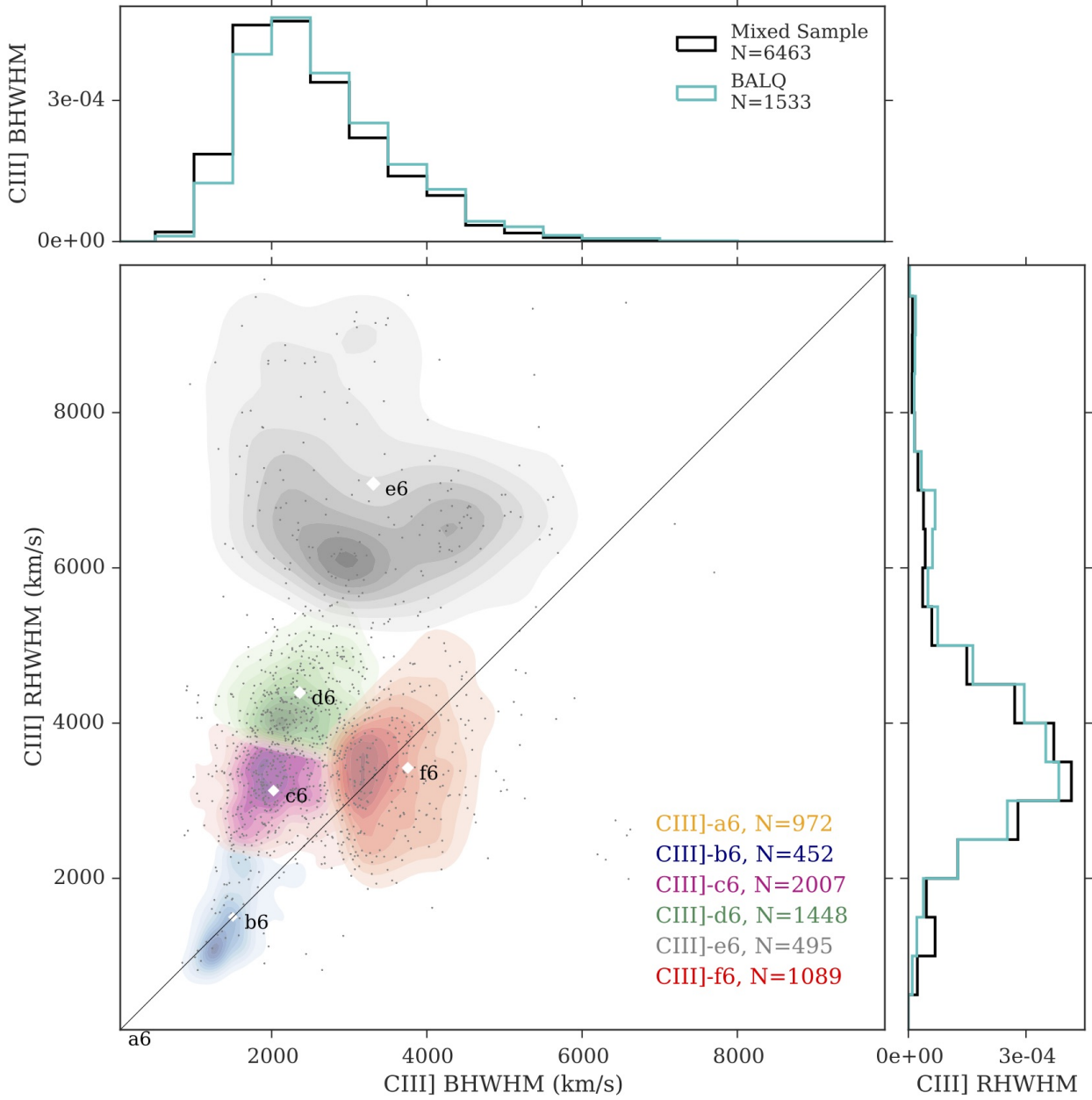


Figure 11. The red- vs. the blue-HWHM of the C III] blend in the mixed sample for the clustering done using $K = 6$. The white diamonds show the projected position of the cluster centroids and their sizes correspond to the average EW of the C III] blend in each cluster. Cluster C III]-a6 (containing objects with failed fitting) is not shown. The BAL quasars (included in this sample) are over-plotted as grey dots and the histograms show the fractional distributions of the BWHM and RWHM in the mixed and the BALQ samples and show no significant difference. The fraction of BAL quasars varies among the clusters: C III]-a6 (33 %), C III]-b6 (11 %), C III]-c6 (19 %), C III]-d6 (24 %), C III]-e6 (33 %), C III]-f6 (23 %). Those differences in the BALQ fractions potentially reflect differences in the hardness of the SED (see §3.4 and Fig. 12).

group of objects (cluster C III]-e5, Figures 9 and 10 in §3.3) with redshifts that are significantly different from those determined using other lines and the PCA redshift determined with the overall fit (see also Table 2). Indeed one of the challenges of working with large datasets is identifying outliers that could potentially be interesting cases or simply inac-

curate measurements that went unnoticed and this rather simple technique offers a straightforward way of identifying such groups.

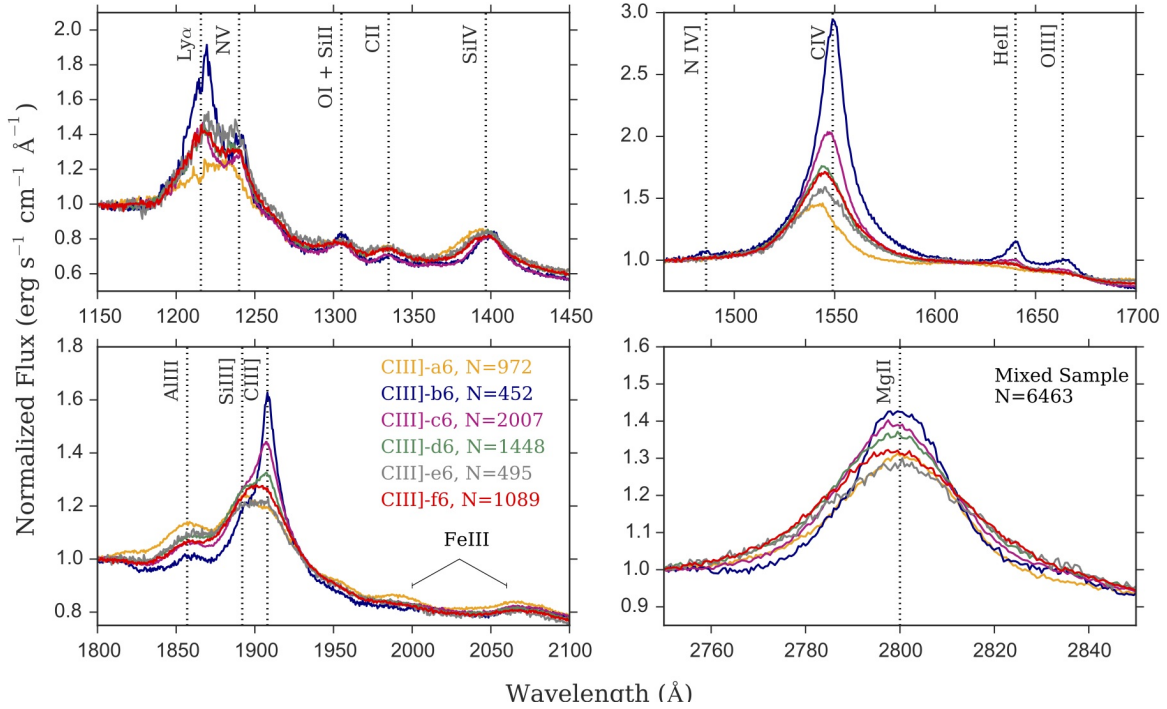


Figure 12. Median composite spectra made from the objects in the C III] clusters shown in Fig. 11 with $K = 6$ in the mixed sample. The profiles have been normalized locally at the starting wavelength of each panel and the numbers in the lower-left panel refer to the number of objects in each composite (cluster). Composite C III]-b6 has prominent narrow features such as C IV, Ly α , and C III] and weak Al III and Si III] (low Si III]/C III] ratio), while composites C III]-e6 and C III]-f6 have weaker C IV that is blueshifted and stronger Si III]/C III] ratios. Composite C III]-a6 is created from the objects with the failed fitting (-1 for the C III] EW, RWHM and BWHM as reported in the catalog). The relatively strong Fe III feature in this cluster might indicate the weak-line quasars comprise a large portion of objects in this clusters.

5 SUMMARY AND CONCLUSION

In this work we explore the use of unsupervised clustering analysis in searching for patterns among a large number of quasar UV spectra in a multidimensional space. The K-means algorithm allows us to find clusters in the EW, BWHM, and RWHM parameter space of three different UV emission lines and line blends: C IV 1450Å, C III] 1903Å, and Mg II 2800Å (see §2.2). We combine objects in each of the clusters we find in median composite spectra to examine the properties of the emission lines in each cluster. As we mention in §2.2, the parameter space we use in this work is expected to be populated continuously and thus K-means is providing us with a way of grouping objects within the continuum that help to reveal the trends from one extreme to the other. Indeed, our composite spectra show that the line properties move gradually between two extremes. The compelling part of the analysis comes from using one line's strength and asymmetries to probe physical properties of the objects (such as the case for the C III] blend probing the hardness of the SED) and seeing the effects of this on other lines not used in the clustering (such as C IV and He II).

We summarize our findings as follows:

- K-means is a simple yet powerful algorithm that, instead of binning data at fixed and rather arbitrary bound-

aries, allows us to more freely explore the structure in a multidimensional parameter space and to find clusters of objects with similar properties in this space.

- When stacked in composite spectra, the clusters we find show some of the well-known trends in quasar UV spectra such as the correlations among the C IV blueshift and its EW and the shape of the ionizing SED probed by the strength of He II and the Si III]/C III] ratio (Figures 7).

- More interestingly, we find this same inverse correlation in the C IV line using emission lines that are seemingly not part of this correlation such as the C III] blend (Fig. 10). Because C III] is generally not contaminated by absorption, it is perhaps more useful than C IV and Mg II for finding like objects in a quasar sample.

- We find that, unlike C IV and C III], the properties of Mg II are not strongly correlated with those of the other lines in the spectra; the width of C IV for example, does not show any clear correlation with that of Mg II. This could lead to potential discrepancy in the determination of black hole masses using broad lines in single epoch spectra (Fig. 5).

- We use this technique to examine the properties of C III] in a mixed sample of BAL and non-BAL quasars and find that the properties of C III] in the mixed sample recovered similar trends to those we find in the main non-BAL quasar sample (Figures 11 and 12).

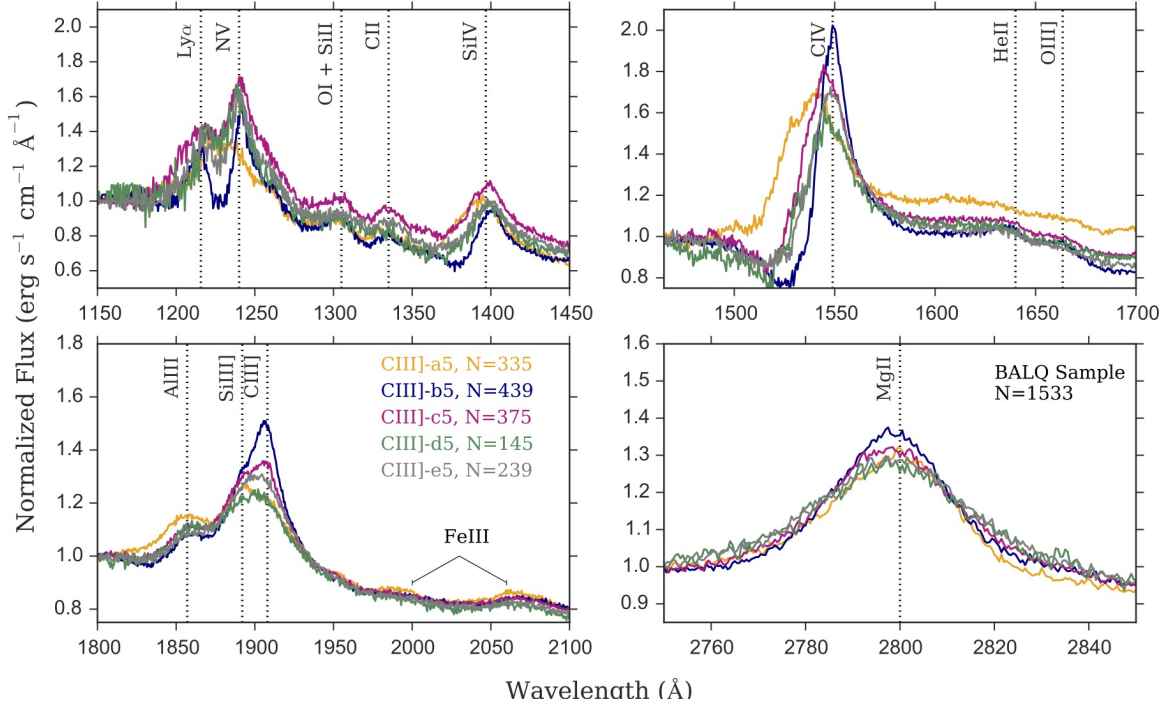


Figure 13. Median composite spectra made from the objects from the BALQ sample in the C III] clusters with $K = 5$. The profiles have been normalized locally at the starting wavelength of each panel and the numbers in the lower-left panel refer to the number of objects in each composite (cluster). Composites made from clusters of the C III] blend properties appear to have different emission and absorption in C IV and Si III], e.g., composite C III]-b5 has the lowest Si III]/C III] ratio and a narrow C IV emission-line that is not blueshifted and a low V_{min} BAL though, while composites C III]-d5 and C III]-e5 have a higher ratio and are both blueshifted in their C IV emission and their BAL troughs have higher V_{min} . Composite C III]-a5 mostly contains objects with -1 measurements for the C III] blend and it shows a weaker C IV emission and a shallow BAL trough with the highest V_{min} . Si IV shows similar trends to those seen in C IV.

- In the mixed sample, we find a higher fraction of BAL quasars in composites with weaker He II and higher Si III]/C III] ratio (softer SEDs), while composites with stronger He II and lower Si III]/C III] ratio (harder SEDs) have a lower fraction of BAL quasars (§3.4, Fig. 11).

- We find that a somewhat more prominent NIV] 1486 Å feature is seen in some of the composites generated from the C III] clusters (composite C III]-a5 in Fig. 10, and composite C III]-a6 in Fig. 12). The stronger He II and higher Si III]/C III] ratio in those composites indicate a hard(er) SED and so the presence of stronger NIV] is not surprising (Karen Leighly, private communication) given that the IP for the NIV] $\lambda 1486$ line is 47.4 eV. Higher nitrogen abundances have been detected before in radio-loud objects (e.g., Jiang et al. 2008) and we find a slightly higher fraction of FIRST-detected objects (12%) in those clusters compare to the ones with weaker nitrogen (7-8%).

- When we apply the clustering to a sample of BAL quasars only using the measurements of C III], we find evidence that the properties of C III] are able to separate objects into groups with different properties of their C IV emission and BAL troughs (Fig. 13). We investigate this result further in Tammour et al. (in prep).

ACKNOWLEDGMENTS

We thank the anonymous referee for helpful comments that improved the manuscript. We are grateful to Patrick Hall and Eric Feigelson for useful discussions at an early stage of this work, Nur Filiz Ak and Nathalie Thibert for comments, and Karen Leighly for thoughtful suggestions that helped improve our manuscript. We also thank Pauline Barnby for inspiring Figures 4, 6, 9, and 11. This work was supported by the Natural Science and Engineering Research Council of Canada, and the Ontario Early Researcher Award Program (A.T., S.C.G.). This research made use of Astropy, a community-developed core Python package for Astronomy (Astropy Collaboration et al. 2013) and Matplotlib: A 2D Graphics Environment (Hunter 2007). Funding for SDSS-III has been provided by the Alfred P. Sloan Foundation, the Participating Institutions, the National Science Foundation, and the U.S. Department of Energy Office of Science. The SDSS-III web site is <http://www.sdss3.org/>.

APPENDIX A: SPECTRA

For completeness, we show here the composite spectra for the rest of the K runs that we did not include in the analysis in §3 for the main sample as well as the mixed and BLAQ samples. The median composite spectra are generated in a similar manner to what is described in §2.3 however the trends discussed in the text are robust to the specific value of K chosen for the clustering.

A1 Main sample -Mg II

A2 Main Sample -C IV

A3 Main sample -C III]

A4 Mixed Sample

A5 BAL Quasars Only

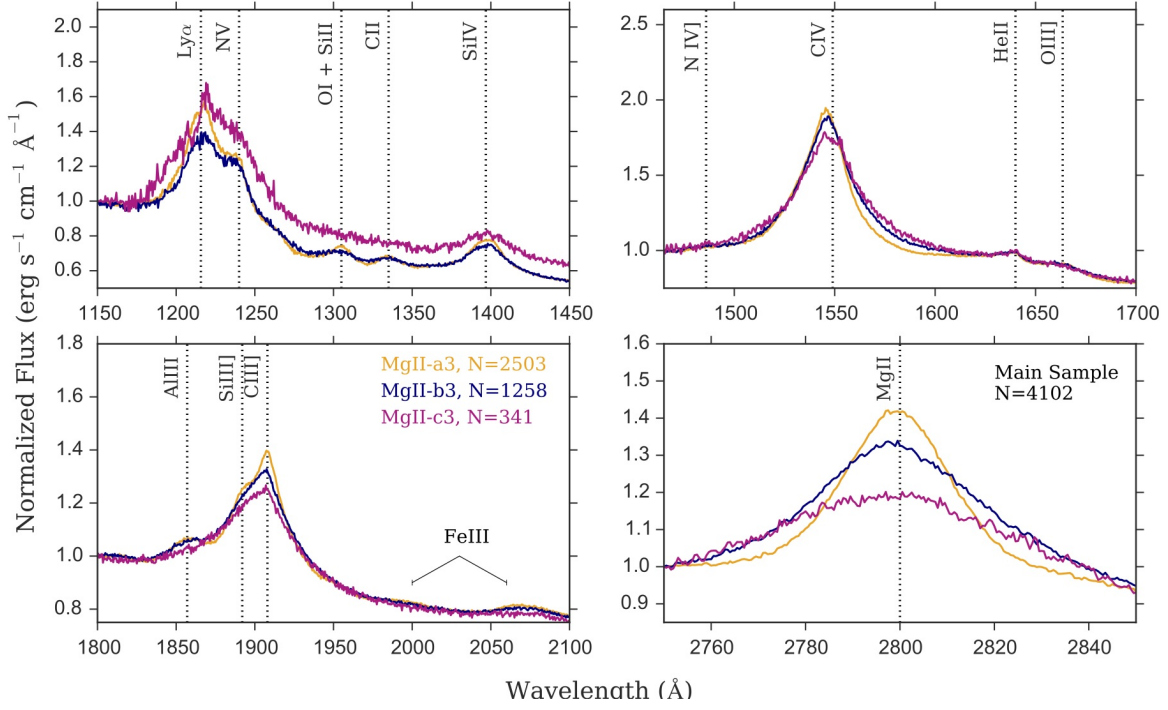


Figure A1. Median composite spectra made from the objects in the Mg II clusters in Fig. 3. Similar to Fig. 5 but for $K = 3$.

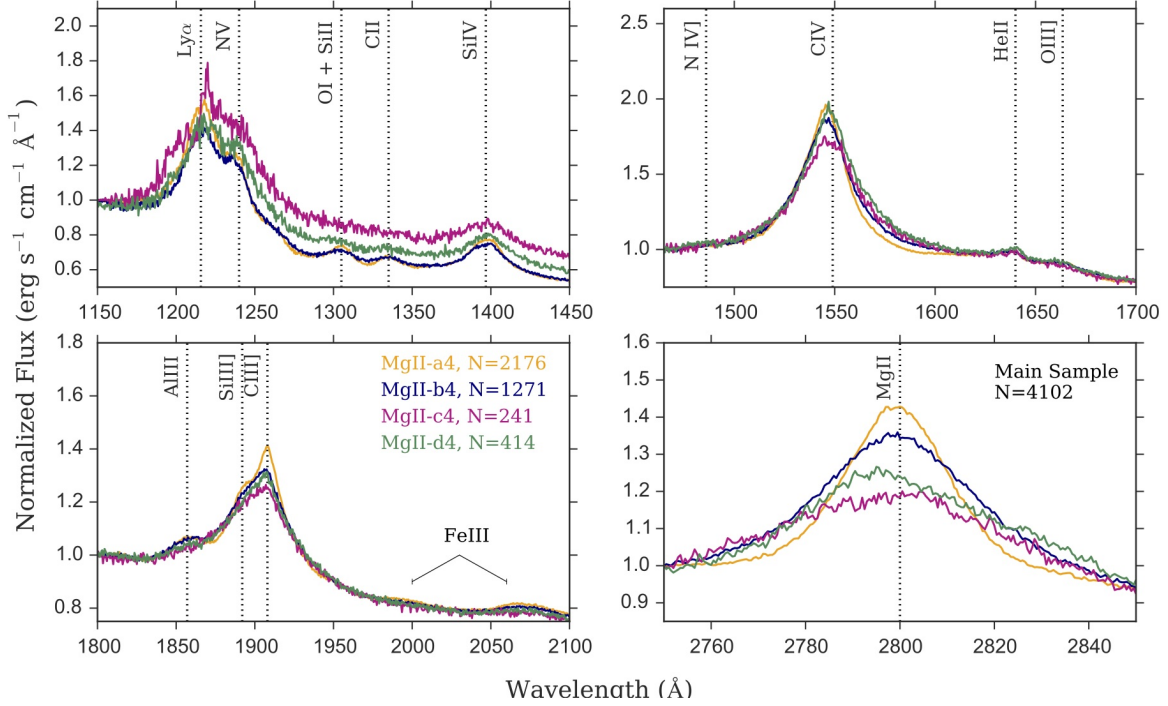


Figure A2. Median composite spectra made from the objects in the Mg II clusters in Fig. 3. Similar to Fig. 5 but for $K = 4$.

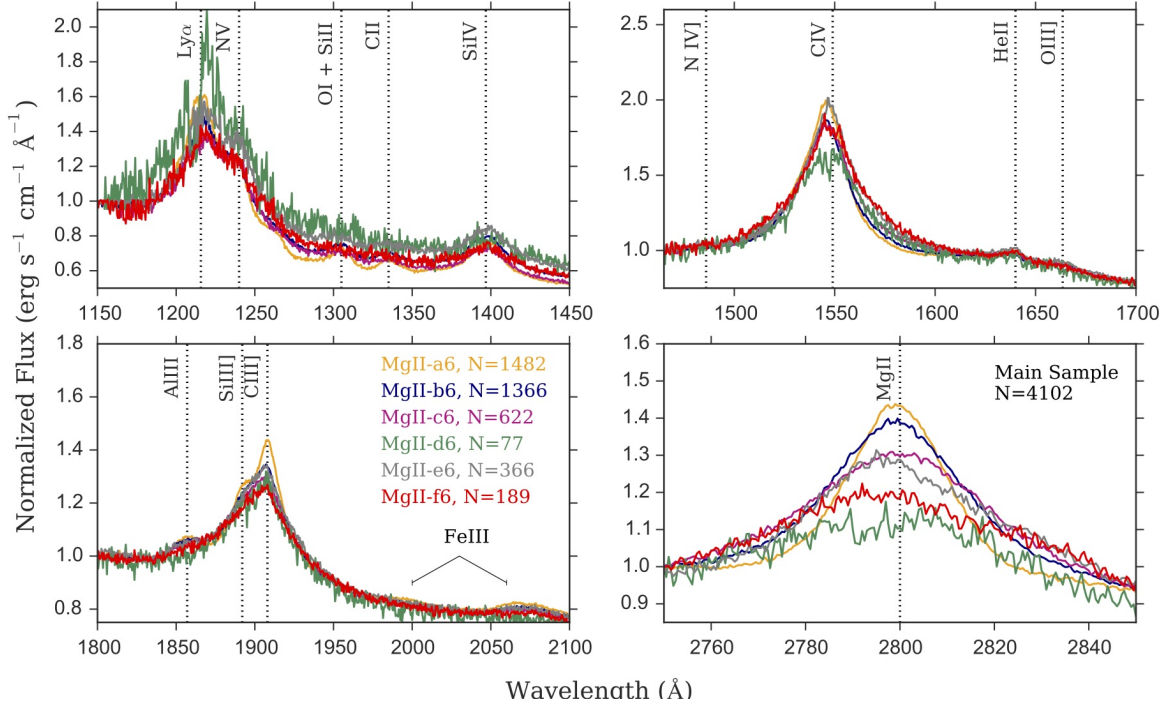


Figure A3. Median composite spectra made from the objects in the Mg II clusters in Fig. 3 Similar to Fig. 5 but for $K = 6$.

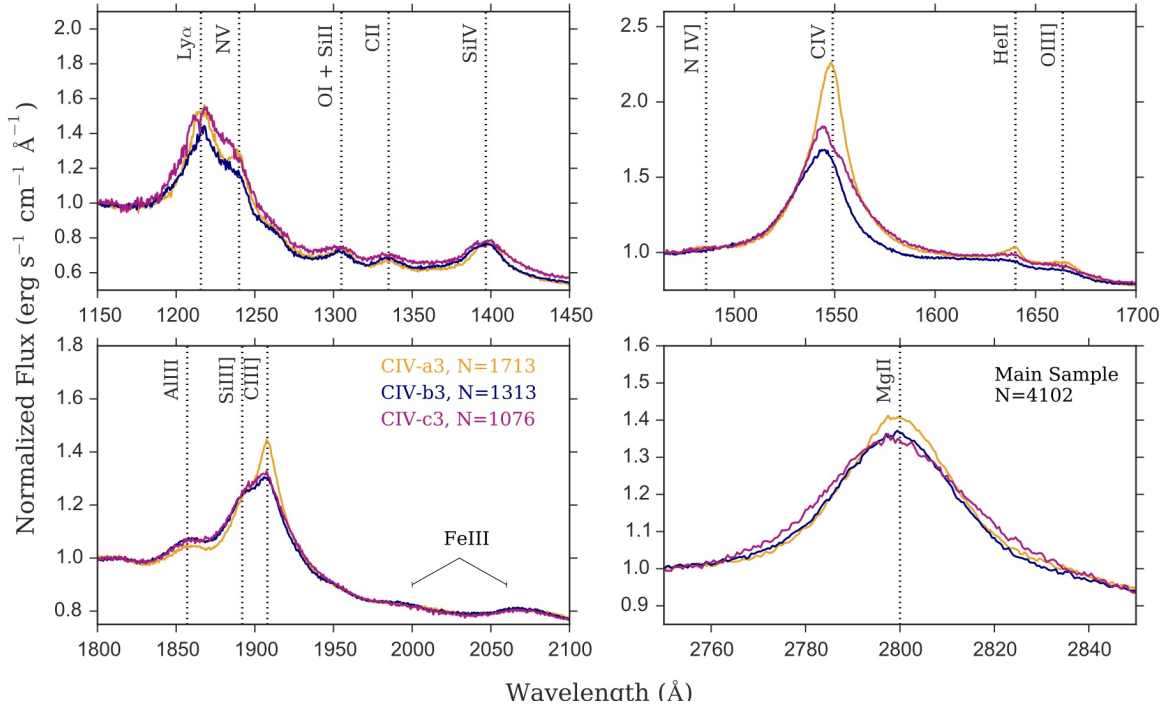


Figure A4. Median composite spectra made from the objects in the C IV clusters. Similar to Fig. 7 but for $K = 3$.

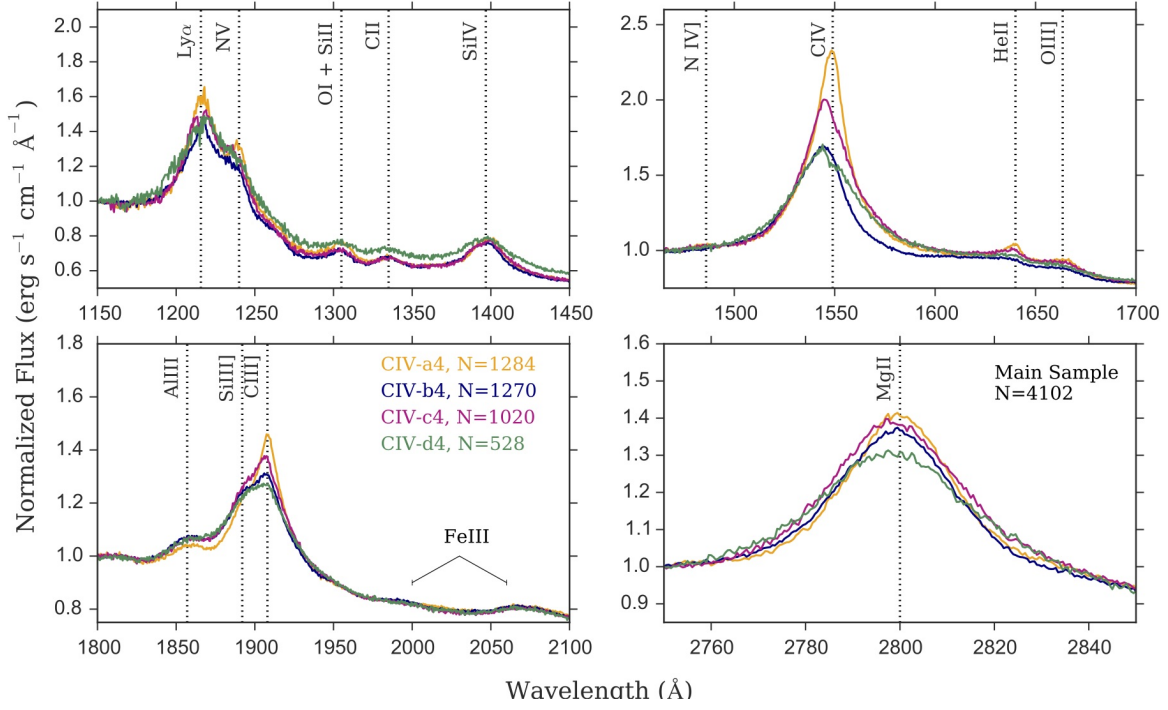


Figure A5. Median composite spectra made from the objects in the C IV clusters. Similar to Fig. 7 but for $K = 4$.

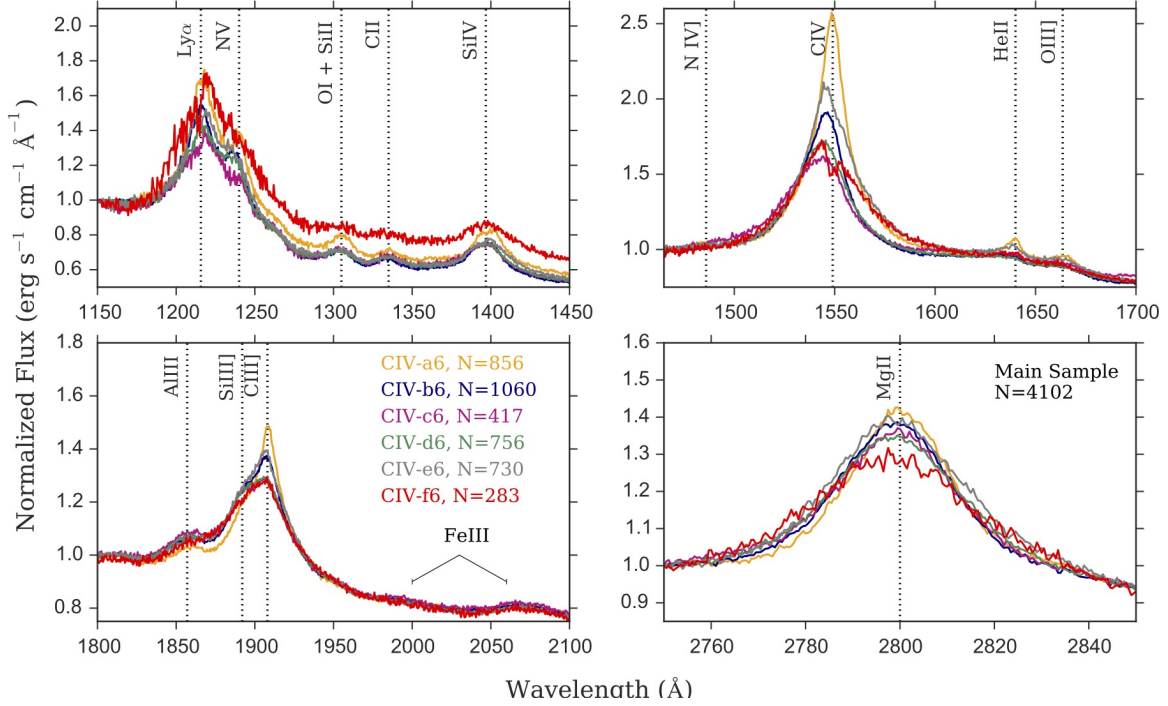


Figure A6. Median composite spectra made from the objects in the C IV clusters. Similar to Fig. 7 but for $K = 6$.

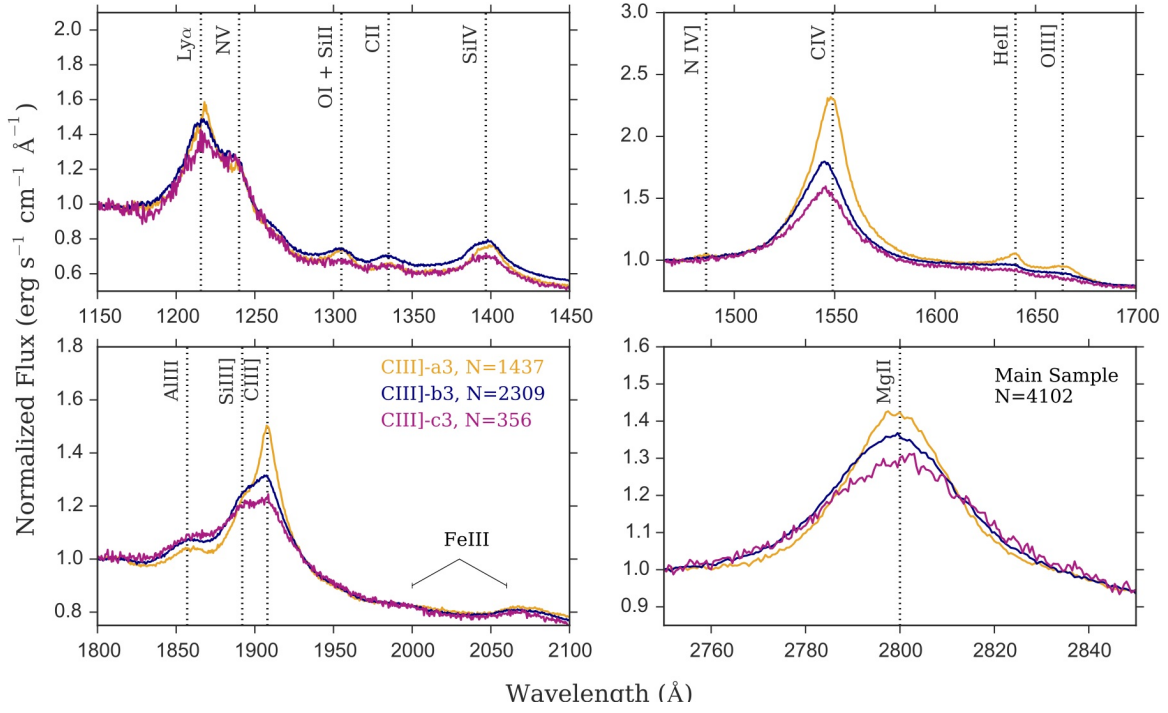


Figure A7. Median composite spectra made from the objects in the C III] clusters. Similar to Fig. 10 but for $K = 3$.

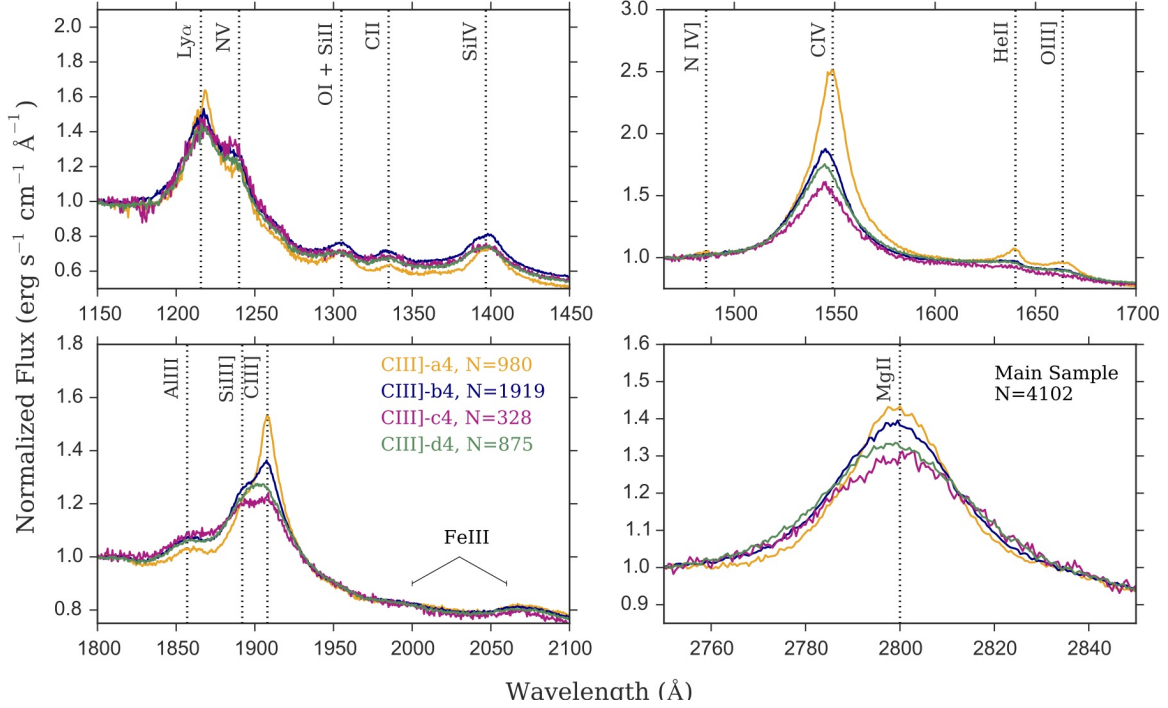


Figure A8. Median composite spectra made from the objects in the C III] clusters. Similar to Fig. 10 but for $K = 4$.

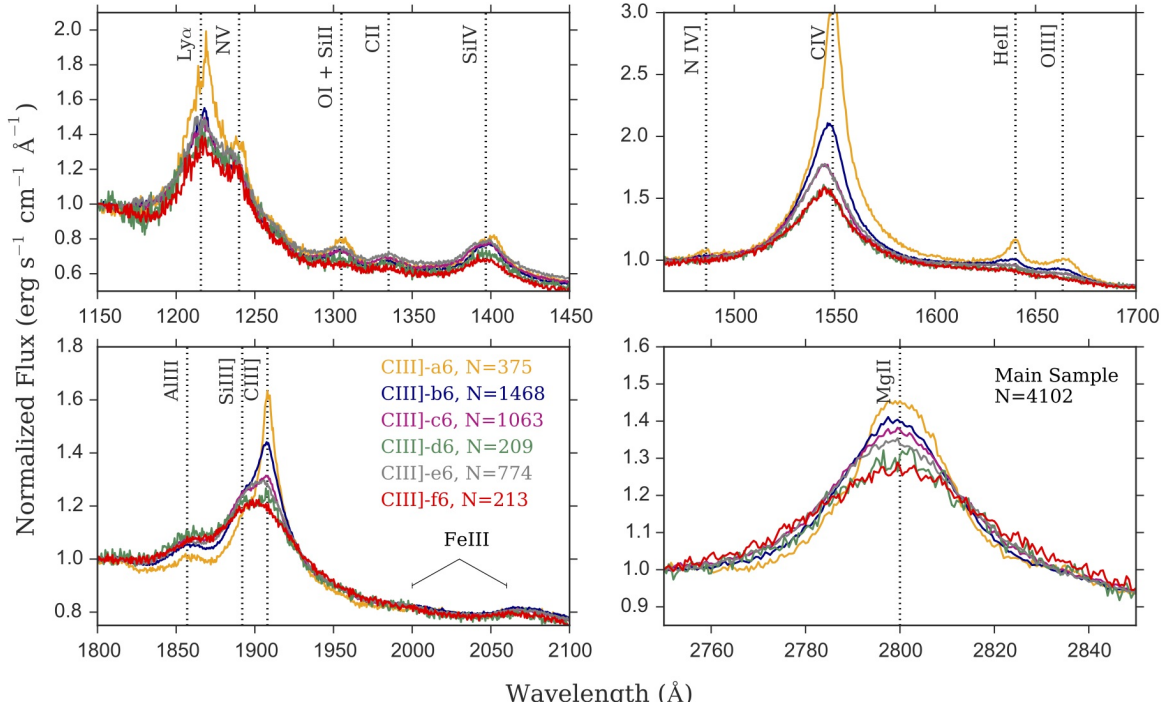


Figure A9. Median composite spectra made from the objects in the C III] clusters. Similar to Fig. 10 but for $K = 6$.

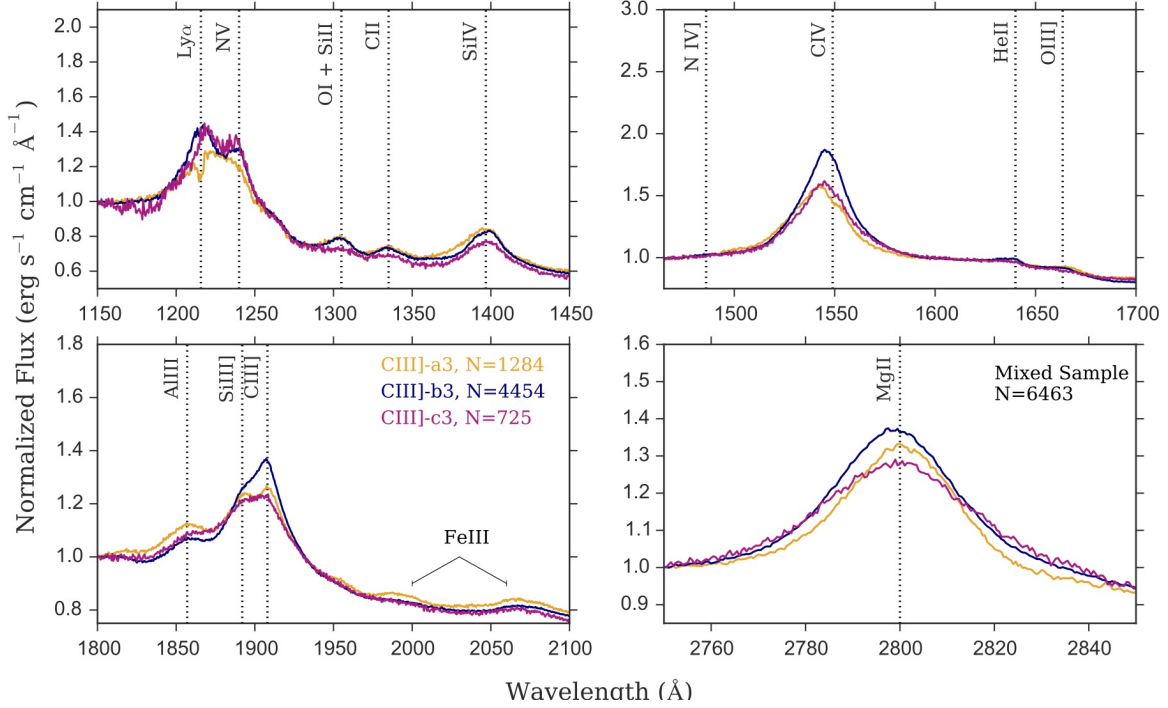


Figure A10. Median composite spectra made from the objects in the C III] clusters in the mixed sample. Similar to Fig. 12 but for $K = 3$.

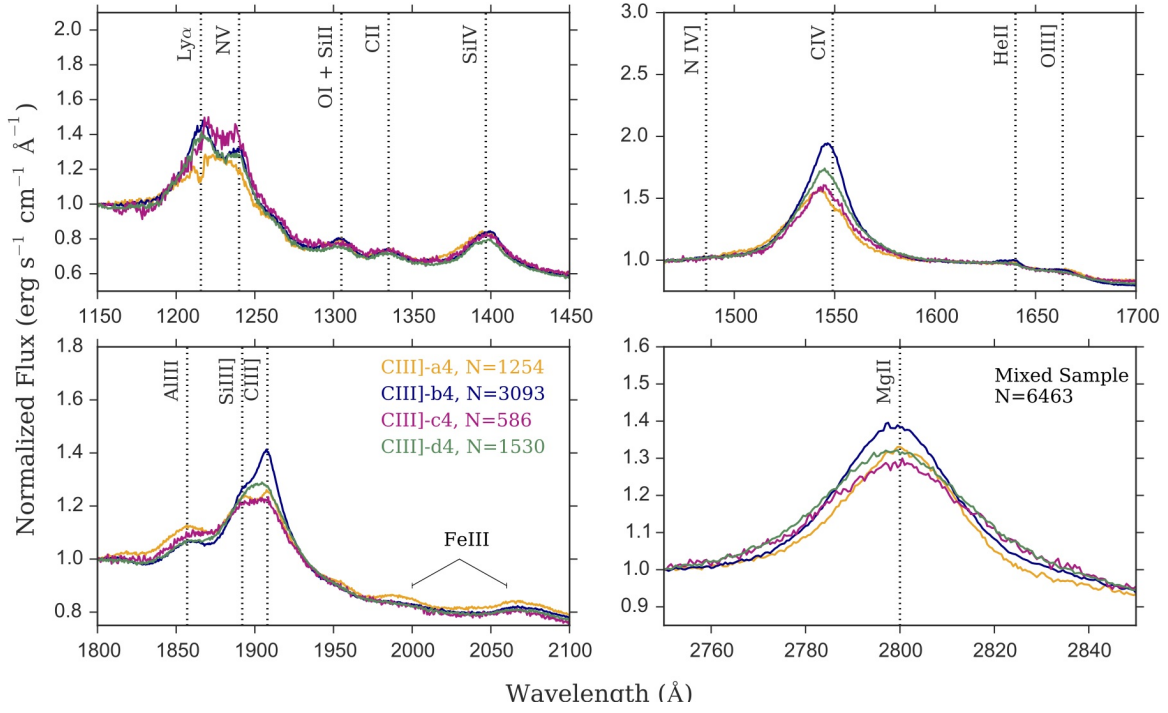


Figure A11. Median composite spectra made from the objects in the C III] clusters in the mixed sample. Similar to Fig. 12 but for $K = 4$.

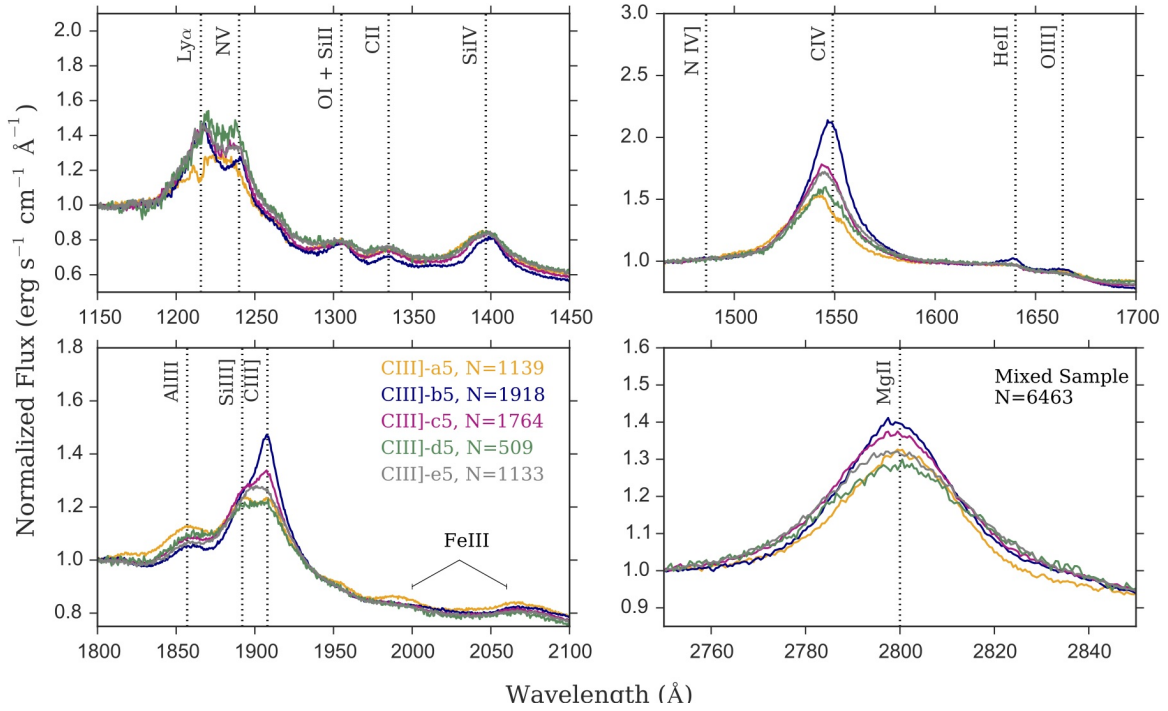


Figure A12. Median composite spectra made from the objects in the C III] clusters in the mixed sample. Similar to Fig. 12 but for $K = 5$.

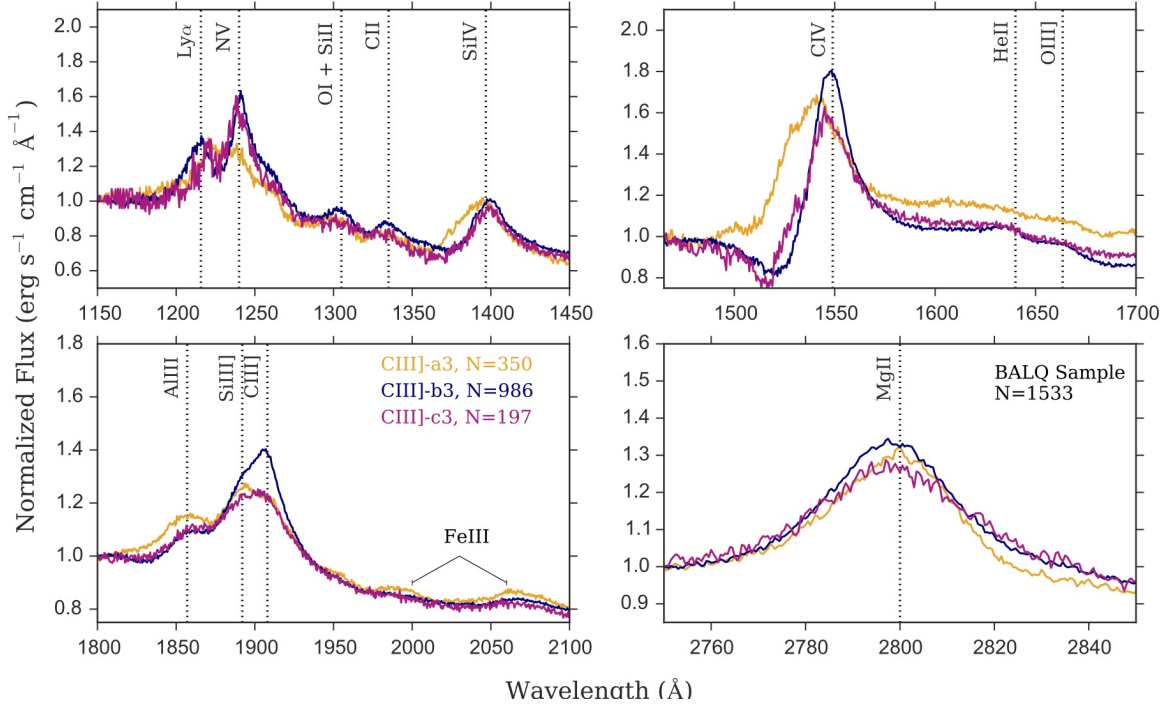


Figure A13. Median composite spectra made from the objects in the C III] clusters in the BALQs sample. Similar to Fig. 13 but for $K = 3$.

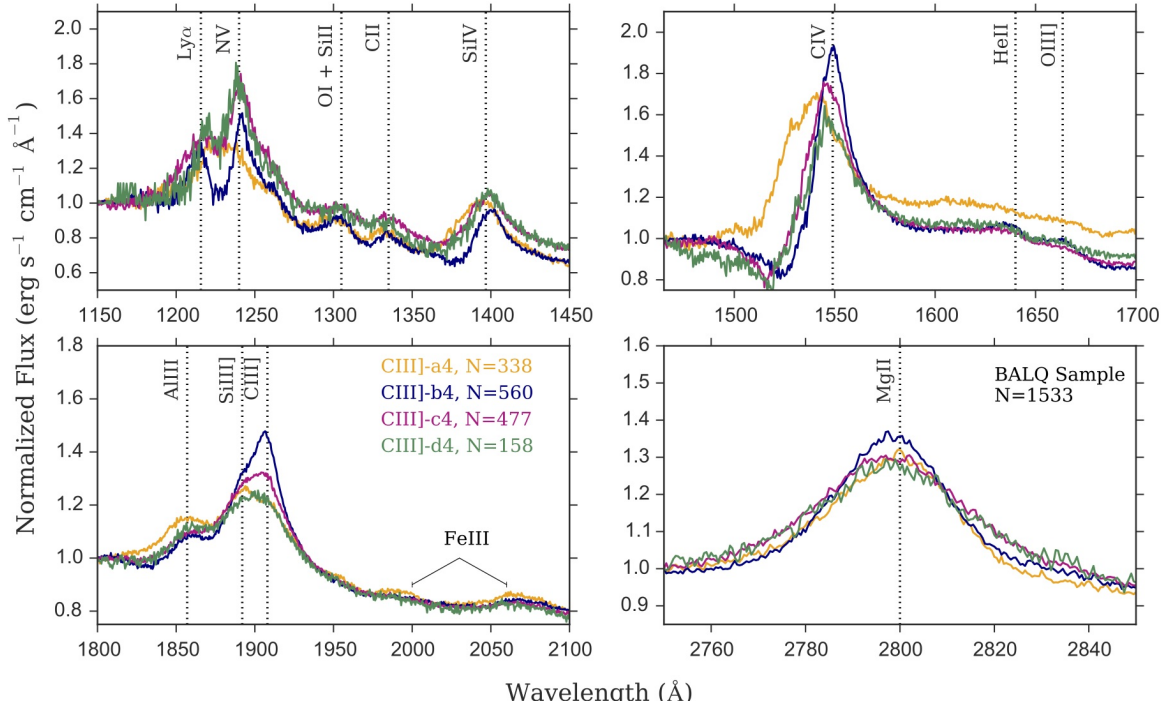


Figure A14. Median composite spectra made from the objects in the C III] clusters in the BALQs sample. Similar to Fig. 13 but for $K = 4$.

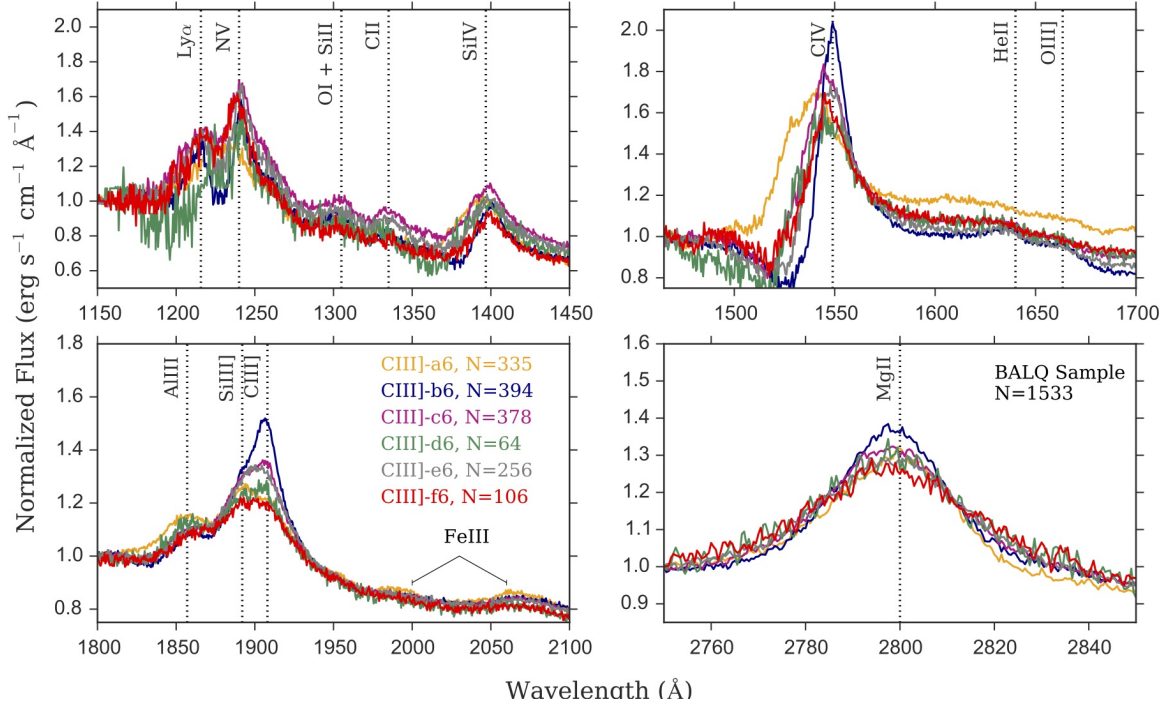


Figure A15. Median composite spectra made from the objects in the C III] clusters in the BALQs sample. Similar to Fig. 13 but for $K = 6$.

APPENDIX B: REPRODUCIBILITY OF CLUSTERS

We tested whether the algorithm is able to find the same unique set of centroids when the clustering is repeated. We find that for the clusters generated in the main sample, the algorithm is finding the same set of clusters even after 50 repeats. We show here one example for $K = 5$.

REFERENCES

- Astropy Collaboration et al., 2013, *A&A*, **558**, A33
- Baldwin J. A., 1977, *ApJ*, **214**, 679
- Baskin A., Laor A., 2005, *MNRAS*, **356**, 1029
- Baskin A., Laor A., Hamann F., 2015, *MNRAS*, **449**, 1593
- Bentz M. C., Hall P. B., Osmer P. S., 2004, *AJ*, **128**, 561
- Bishop C., 2009, *Pattern Recognition and Machine Learning*. Springer
- Boroson T. A., 2002, *ApJ*, **565**, 78
- Boroson T. A., Green R. F., 1992, *ApJS*, **80**, 109
- Casebeer D. A., Leighly K. M., Baron E., 2006, *ApJ*, **637**, 157
- Collin-Souffrin S., Dyson J. E., McDowell J. C., Perry J. J., 1988, *MNRAS*, **232**, 539
- Croom S. M., et al., 2002, *MNRAS*, **337**, 275
- Gallagher S. C., Brandt W. N., Chartas G., Priddey R., Garmire G. P., Sambruna R. M., 2006, *ApJ*, **644**, 709
- Han J., Kamber M., Pei J., 2012, *Data Mining: Concepts and Techniques*. Elsevier
- Hastie T., Tibshirani R., Friedman J., 2009, *The Elements of Statistical Learning*. Springer
- Hewett P. C., Foltz C. B., 2003, *AJ*, **125**, 1784
- Hewett P. C., Wild V., 2010, *MNRAS*, **405**, 2302
- Hill A. R., Gallagher S. C., Deo R. P., Peeters E., Richards G. T., 2014, *MNRAS*, **438**, 2317
- Hunter J. D., 2007, *Computing In Science and Engineering*, **9**, 90
- Ivezić Z., Connolly A. J., VanderPlas J., Gray A., 2014, *Statistics, Data Mining, and Machine Learning in Astronomy: A Practical Python Guide for the Analysis of Survey Data*. Princeton University Press
- Jiang L., Fan X., Vestergaard M., 2008, *ApJ*, **679**, 962
- Kruczek N. E., et al., 2011, *AJ*, **142**, 130
- Leighly K. M., 2004, *ApJ*, **611**, 125
- Leighly K. M., Moore J. R., 2004, *ApJ*, **611**, 107
- Leighly K. M., Halpern J. P., Jenkins E. B., Casebeer D., 2007, *ApJS*, **173**, 1
- Luo B., et al., 2015, *ApJ*, **805**, 122
- Murray N., Chiang J., 1998, *ApJ*, **494**, 125
- Pâris I., et al., 2012, *A&A*, **548**, A66
- Pâris I., et al., 2014, *A&A*, **563**, A54
- Pedregosa F., et al., 2011, *Journal of Machine Learning Research*, **12**, 2825
- Proga D., Stone J. M., Kallman T. R., 2000, *ApJ*, **543**, 686
- Richards G. T., Vanden Berk D. E., Reichard T. A., Hall P. B., Schneider D. P., SubbaRao M., Thakar A. R., York D. G., 2002, *AJ*, **124**, 1
- Richards G. T., et al., 2011, *AJ*, **141**, 167
- Shakura N. I., Sunyaev R. A., 1973, *A&A*, **24**, 337
- Shemmer O., Lieber S., 2015, *ApJ*, **805**, 124
- Shen Y., Ho L. C., 2014, *Nature*, **513**, 210
- Shen Y., Liu X., 2012, *ApJ*, **753**, 125
- Spergel D. N., et al., 2003, *ApJS*, **148**, 175
- Sulentic J. W., Zwitter T., Marziani P., Dultzin-Hacyan D., 2000, *ApJ*, **536**, L5
- Sulentic J. W., Bachev R., Marziani P., Negrete C. A., Dultzin D., 2007, *ApJ*, **666**, 757
- Tammour A., Gallagher S. C., Richards G., 2015, *MNRAS*, **448**, 3354
- Vanden Berk D. E., et al., 2001, *AJ*, **122**, 549
- Weymann R. J., Morris S. L., Foltz C. B., Hewett P. C., 1991, *ApJ*, **373**, 23
- Yip C. W., et al., 2004, *AJ*, **128**, 2603

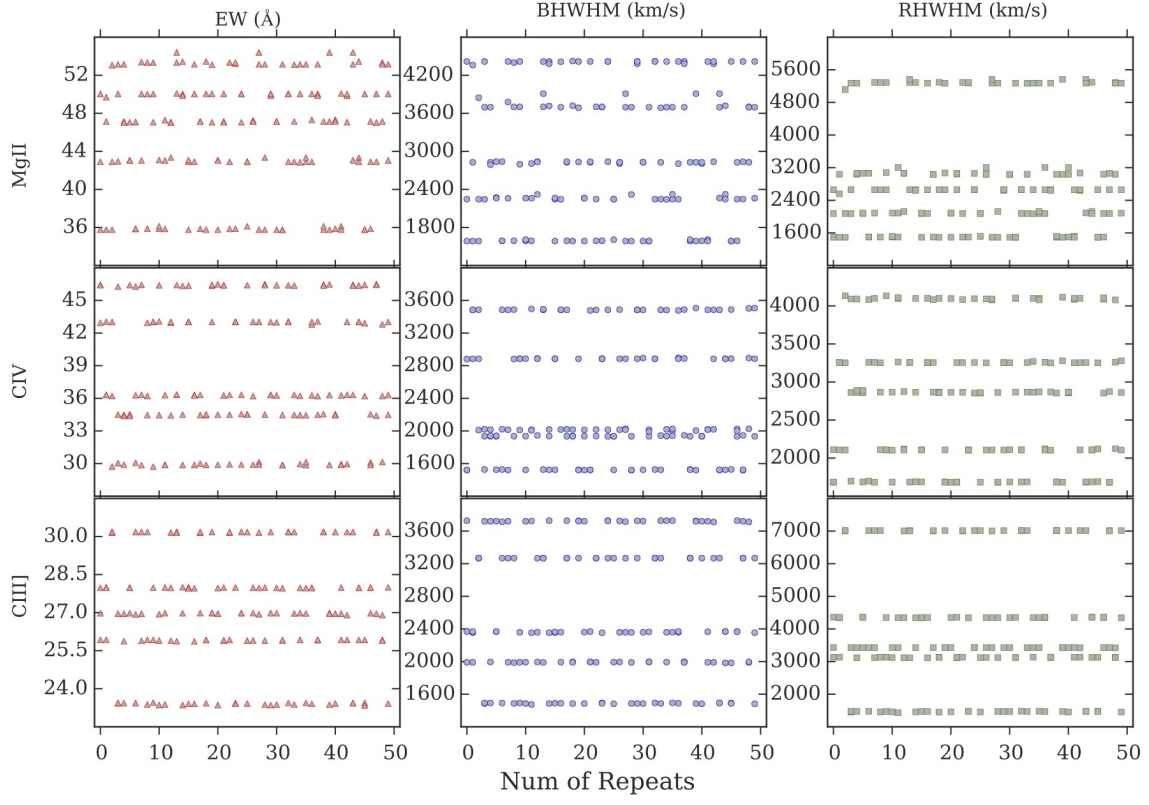


Figure B1. The centroids for three different clustering attempts one for each of the lines (Mg II, C IV, and C III) in the main sample repeated 50 times using $K = 5$. This shows that K-means is finding clusters that are consistent and that the clusters we use in our analysis reflect meaningful groupings of like objects in the parameter space of EW, RWHM and BWHM.